

# A formal introduction to RL theory

Sebastian Scherer

November 12, 2019

## Contents

1	Why another introduction?	2
2	Some notation	2
3	Value function and action-value function	3
3.1	Existential crisis . . . . .	4
3.2	Relationships . . . . .	10
3.3	Comparing and improving . . . . .	13
4	Optimal policies and how to find them	20
4.1	Existential crisis . . . . .	20
4.2	Behavioural issues . . . . .	21
4.3	The Bellman operator . . . . .	24
4.4	Building optimal policies . . . . .	26

## 1 Why another introduction?

The goal of this work is to provide a mathematically rigorous introduction to RL theory based on the excellent book "Introduction to reinforcement learning" by Sutton and Barto (first edition). While I greatly enjoyed reading this book and appreciated its focused approach on developing an intuition for the Q- and V-functions, the algorithms and the general probabilistic framework introduced in the early chapters, I couldn't help but stumble at some points wondering how exactly a particular claim was justified. Queries on stack exchange as well as the various alternative resources applying even less rigour and, often times, introducing additional confusing notation, motivated me to try and remedy this myself. I therefore set out to try and formalize the theory presented, at least for the finite Markov Decisions Processes treated in the book, so that it may help let my inner mathematician sleep at night, as well as, and this is my sincere hope, provide a rigorous and helpful introduction for all those who are not only interested in the intuition but also appreciate a firm foundation on which to place it. The following manuscript can be used as an explanatory guide to the concepts presented in the book, or can be independently used as a rigorous introduction to value function theory in its own right.

## 2 Some notation

Like the reference book, we consider finite state, finite action markov decision processes ("finite MDPs"). As such, we denote by  $S$  the set of states achievable for a given finite MDP, and by  $A$  the set of executable actions  $a$ . We do not restrict ourselves to deterministic policies, and therefore treat a policy  $\pi$  as a conditional probability distribution over the executable action set  $A$ , conditioned on a given current state from  $S$ . In other words,

$$\begin{aligned} \pi &: A \times S \rightarrow [0, 1] \\ (a, s) &\mapsto \pi(a, s) \end{aligned} \tag{1}$$

where  $\pi(a, s) = Pr_{\pi}(a|s)$  denotes the probability of choosing action  $a \in A$  when in state  $s \in S$  while acting according to policy  $\pi$ .

We encode our knowledge about the (reactionary) nature of our environment via the transition probabilities

$$P_{s,s'}^a = Pr(s_{t+1} = s' | s_t = s, a_t = a) \tag{2}$$

where  $s, s' \in S$  and  $a \in A$ , denoting the probability of ending up in state  $s'$  at  $t + 1$  when coming from state  $s$  at  $t$  by executing  $a$ , and

$$R_{s,s'}^a = \mathbb{E}[r_t | s_t = s, s_{t+1} = s', a_t = a] \tag{3}$$

denoting the expected reward at time  $t$  due to ending up in state  $s'$  at  $t + 1$  when coming from state  $s$  at  $t$  by executing  $a$ . For the remainder of this script, we will assume a uniform bound on rewards throughout all time steps, i.e. that

$$-M \leq r_t \leq M \quad (4)$$

for some  $M \in \mathbb{R}$  and each  $t \in 0, 1, 2, \dots$ .

Note that, in our notation,  $s_t$  and  $a_t$  denote the state and action at time  $t$  respectively, and thus  $r_t$  - NOT  $r_{t+1}$  as in the book - denotes the reward obtained AFTER being in  $s_t$  and executing  $a_t$ , thereby resulting in some (possibly the same) state  $s_{t+1}$ .

We use the same symbol  $\gamma \in (0, 1)$  to denote the reward discount factor.

Finally, the most difficult notation to right *and* consistent: expected values. We will use slightly different notations to indicate the various different underlying distributions that govern the behaviour of the random variables involved, and w.r.t which the expected value needs to be viewed.

If we are dealing with an implicit sequence of actions chosen according to one policy like

$$s_t \xrightarrow{\pi} a_t \xrightarrow{P_{s_t, \cdot}^{a_t}} s_{t+1} \xrightarrow{\pi} a_{t+1} \xrightarrow{P_{s_{t+1}, \cdot}^{a_{t+1}}} s_{t+2} \xrightarrow{\pi} \dots, \quad (5)$$

we will express this by writing  $\mathbb{E}_\pi[\cdot]$ . The contribution of the environment's state distribution  $P_{s_t, \cdot}^{a_t}$  is implicit since we usually deal with one finite MDP at a time, thereby keeping this particular distribution constant throughout all of the proofs. An example of this is

$$\mathbb{E}_\pi[r_t | s_t = s],$$

which implies action  $a_t$  was taken according to  $\pi$  having started at state  $s$  at time  $t$ .

Therefore, if the random variable in question is only dependent on the environment's distribution, such as in the expression

$$\mathbb{E}[r_t + \gamma f(s_{t+1}) | s_t = s, a_t = a]$$

(where  $f$  is some deterministic function) we omit any index. Note that in this case both the immediate reward  $r_t$  as well as the next time step's state  $s_{t+1}$  are entirely dependent on the environment parameters  $R_{s, s'}^a$  and  $P_{s, s'}^a$ , since we fixed the action  $a_t$ , thereby cutting any potential policy out of the loop.

In some cases, we will need to indicate the involved random variables' distributions explicitly and individually. In those cases, we will make clear what exactly we mean.

### 3 Value function and action-value function

As Sutton and Barto point out, the standard approach to the analysis of optimal behaviour w.r.t a given MDP is by closely examining the value function  $V^\pi$  and action value function  $Q^\pi$  associated to a given policy  $\pi$ . They are an essential tool for quantitative analysis of policy driven behaviour, and thus, unsurprisingly, we will make heavy use of them throughout this guide. For a given policy  $\pi$  on a finite MDP, we call

$$\begin{aligned} V^\pi : S &\rightarrow \mathbb{R} \\ s &\mapsto \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s] \end{aligned} \quad (6)$$

the *value function* of  $\pi$ , and

$$\begin{aligned} Q^\pi : A \times S &\rightarrow \mathbb{R} \\ (a, s) &\mapsto \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s, a_t = a] \end{aligned} \quad (7)$$

the *action-value* function of  $\pi$ . The general idea of  $V^\pi$  is the quantification of the value a certain state  $s$  possesses under  $\pi$ , by assigning it the expected cumulative reward obtained when starting in that state  $s$  and then acting (i.e. choosing and executing actions) according to  $\pi$ .  $Q^\pi$  does very much the same thing, except for pairs of states and actions  $(a, s)$ , assigning expected cumulative rewards when starting from  $s$  via action  $a$ , and only *then* acting according to  $\pi$ .

### 3.1 Existential crisis

In line with our rigorous approach so far, we dedicate this subsection to proving the well definedness of the above quantities. That is, we get assurances for:

- the well definedness of the random variable defined by the series  $G_t := \sum_{k=0}^{\infty} \gamma^k r_{t+k}$
- the well definedness of its expected value, i.e. the well definedness of  $V^\pi$
- the well definedness of  $Q^\pi$ .
- some handy derivative results from the above three

We start our quest with guaranteeing the well definedness of  $G_t$ . As a random variable that is the discounted infinite sum of random variables  $r_{t+k}$ ,  $k = 0, 1, \dots$ , it maps sequences  $\omega = (s_t, a_t, s_{t+1}, a_{t+1}, \dots)$ , with  $s_t, s_{t+1}, \dots \in S$ ,  $a_t, a_{t+1}, \dots \in A$  onto the set of real numbers  $\mathbb{R}$ . In other words, by denoting  $\Omega$  as the set of all such sequences  $\omega$ , we consider  $G_t$  to be

$$\begin{aligned} G_t : \Omega &\rightarrow \mathbb{R} \\ \omega &\mapsto \sum_{k=0}^{\infty} \gamma^k r_{t+k}(\omega) \end{aligned} \quad (8)$$

where  $\sum_{k=0}^{\infty} \gamma^k r_{t+k}(\omega) \stackrel{!}{\in} \mathbb{R}$  remains to be shown to guarantee well definedness of this quantity as a (random variable) mapping. We will do so in the following

**Proposition 1.** (*Cumulative rewards exist*) *The random variable  $G_t$  as defined in 8 is well defined. In other words, the series  $G_t(\omega) = \sum_{k=0}^{\infty} \gamma^k r_{t+k}(\omega)$  is convergent for all  $\omega \in \Omega$ .*

*Proof.* We will show that for any trajectory  $\omega \in \Omega$  the sequence  $G_t^n(\omega) := \sum_{k=0}^n \gamma^k r_{t+k}(\omega) \in \mathbb{R}$  is a Cauchy sequence. Since  $\mathbb{R}$  is complete, this will imply convergence of  $G_t^n(\omega)$  in  $\mathbb{R}$ .

Let w.l.o.g  $n > m$  to see that

$$|G_t^n(\omega) - G_t^m(\omega)| = \left| \sum_{k=m}^n \gamma^k r_{t+k}(\omega) \right| \stackrel{4}{\leq} M \sum_{k=m}^n \gamma^k. \quad (9)$$

Since  $s_n := \sum_{k=0}^n \gamma^k \xrightarrow{n \rightarrow \infty} \frac{1}{1-\gamma}$  and therefore cauchy, there exists for every  $\hat{\epsilon}$  an  $N_{\hat{\epsilon}}$  such that, if,  $n, m \geq N_{\hat{\epsilon}}$ ,  $\sum_{k=m}^n \gamma^k = |s_n - s_m| < \hat{\epsilon}$ . If now  $\epsilon > 0$  is arbitrary but fixed, choose  $\hat{\epsilon} = \frac{\epsilon}{M}$ . Thus, for any  $n, m > N_{\hat{\epsilon}}$ , we have

$$|G_t^n(\omega) - G_t^m(\omega)| \stackrel{9}{\leq} M|s_n - s_m| < M \frac{\epsilon}{M} = \epsilon. \quad (10)$$

Thus for any  $\omega \in \Omega$ , the sequence  $G_t^n(\omega)$  is cauchy, and thus convergent on  $\mathbb{R}$ . This proves the claim.  $\square$

Having established  $G_t(\omega) \in \mathbb{R}$ , we can immediately see that it must be uniformly bounded on  $\Omega$ .

**Corollary 1.** (*Cumulative rewards are bounded*) *The cumulative reward  $G_t$  as defined in 8 is uniformly bounded on  $\Omega$ . In particular,*

$$|G_t(\omega)| < \frac{M}{1-\gamma} \quad (11)$$

for all  $\omega \in \Omega$ .

*Proof.* Thanks to Prop 1 we can justify writing

$$|G_t(\omega)| = \left| \sum_{k=0}^{\infty} \gamma^k r_{t+k}(\omega) \right| \leq \sum_{k=0}^{\infty} \gamma^k |r_{t+k}(\omega)| \stackrel{4}{=} M \sum_{k=0}^{\infty} \gamma^k = \frac{M}{1-\gamma}. \quad (12)$$

$\square$

Now that the infinite discounted reward is a well defined random variable, we can sensibly pose the question regarding its expected value. Before we do so, however, let us have a look at the trajectory space  $\Omega$  and consider why it is a countable, discrete set.

Let us denote by  $c_S = |S|$ ,  $c_A = |A|$  the cardinality, i.e. size, of the state space  $S$  and action space  $A$ , respectively. We can then order all possible trajectories in

$$\Omega = \{\omega = (s_1, a_2, s_2, a_2, \dots) | s_i \in S, a_i \in A, i = 1, 2, \dots\}. \quad (13)$$

To see such an ordering, refer to the appendix (? - still have to write this part). In particular, whenever we write  $\omega_i$  for  $i = 1, 2, \dots$  or similar we mean the ordered  $\omega \in \Omega$ .

We now establish the existence of the expected value of  $G_t$ .

**Lemma 1.** (*Expected cumulative rewards exist, too*) *Let  $P_{\Omega}$  be any trajectory distribution on the trajectory space  $\Omega$ . Then the expected value of the cumulative reward  $G_t$  under that distribution is well defined and exists. In other words  $\mathbb{E}_{\omega \sim P_{\Omega}}[G_t] < \infty$ .*

Note that this is a pretty general statement about sufficient conditions for the existence of  $G_t$ 's expected value. It includes, for example, the cases where we act according to some policy  $\pi$ , starting from some random state  $s$  distributed according to some random distribution  $P_S$  on  $S$ . In particular, if that distribution is deterministic, i.e. if

$$P_{S,s'}(s) = \begin{cases} 1 & \text{for } s = s' \\ 0 & \text{else} \end{cases} \quad (14)$$

we will have established the well-definedness of  $V_\pi$  as defined in Definition 6, since that is what the value function, according to 6, measures. To see this, note that in such a case,  $P_\Omega = P_{\Omega|s'}^\pi = P^\pi \circ P_{S,s'}$  is the distribution resulting from sampling a starting state  $s$  according to  $P_{S,s'}$ , and then, given such a starting state, sampling a trajectory  $\omega$  according to  $\pi$ . In such a case we have

$$V^\pi(s') \stackrel{Def6}{=} \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s \right] = \mathbb{E}_{P_\Omega}[G_t] \stackrel{Lem1}{<} \infty \quad (15)$$

The analogous argument guaranteeing the well definedness of our action-value function  $Q^\pi$  as defined in Definition 7, uses a deterministic distribution  $P_{S \times A, (s', a')}$  such that, for a fixed but arbitrary pair  $(s', a') \in S \times A$ ,

$$P_{S \times A, (s', a')}(s, a) = \begin{cases} 1 & \text{for } (s, a) = (s', a') \\ 0 & \text{else.} \end{cases} \quad (16)$$

This leads to a composite trajectory distribution  $P_\Omega = P_{\Omega|s', a'}^\pi = P^\pi \circ P_{S, (s', a')}$  resulting from sampling a starting state *and* action  $(s, a)$  according to  $P_{S, (s', a')}$ , and then, given such a starting state-action pair, sampling a trajectory  $\omega$  according to  $\pi$ . This yields

$$Q^\pi(s', a') \stackrel{Def7}{=} \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s', a_t = a' \right] = \mathbb{E}_{P_\Omega}[G_t] \stackrel{Lem1}{<} \infty. \quad (17)$$

For easier reference later on, let's summarize these arguments in the form of the following

**Corollary 2.** *((Action-)Value functions exist) Let  $\pi$  be probabilistic or deterministic policy on a finite MDP. Then the functions  $V^\pi$  as defined in 6 and  $Q^\pi$  as defined in 7 are well defined.*

*Proof.* For a proof of both claims, see the argumentats presented in the paragraph preceding this corollary. In particular, for the claim about  $V^\pi$  see the derivation of equations 14 and 15. Similarly, the well definedness of  $Q^\pi$  is justified by the arguments formalized in equations 16 and 17.  $\square$

Having established the usefulness of this Lemma, let us finally prove it.

*Proof.* (of Lemma 1) Given some trajectory distribution  $P_\Omega$  on the (unordered) set  $\Omega$ , we need to show convergence of the series

$$\mathbb{E}_{\omega \sim P_\Omega}[G_t] = \sum_{\omega \in \Omega} P_\Omega(\omega) G_t(\omega) < \infty \quad (18)$$

regardless of the particular ordering of the trajectories  $\omega$ . We do this by showing the convergence of this series' absolute series counterpart, namely

$$\sum_{\omega \in \Omega} |P_{\Omega}(\omega)G_t(\omega)| \stackrel{P_{\Omega}(\omega) \geq 0}{=} \sum_{\omega \in \Omega} P_{\Omega}(\omega)|G_t(\omega)| = \sum_{i=1}^{\infty} P_{\Omega}(\omega_i)|G_t(\omega_i)| < \infty \quad (19)$$

for the *one particular ordering* of  $\Omega$  referred to earlier. The convergence of the absolute series in 19 for one possible ordering of  $\omega \in \Omega$  then guarantees the convergence of the series in 18 for *any* ordering of  $\omega \in \Omega$  to the same limit according to the "Famous Reordering Series Convergence Theorem". To show that the absolute series converges, we show that the sequence

$$EG_{abs}^n := \sum_{i=1}^n |G_t(\omega_i)|P_{\Omega}(\omega_i) \stackrel{Prop1}{<} \infty \quad (20)$$

induced by its partial sum is a Cauchy sequence. W.l.o.g, let  $n > m$  to see that

$$\begin{aligned} |EG_{abs}^n - EG_{abs}^m| &= \sum_{i=m}^n |G_t(\omega_i)|P_{\Omega}(\omega_i) \\ &\stackrel{Cor1}{\leq} \frac{M}{1-\gamma} \sum_{i=m}^n P_{\Omega}(\omega_i). \end{aligned} \quad (21)$$

Since  $\sum_{i=0}^{\infty} P_{\Omega}(\omega_i) = 1$ , the sequence of partial sums

$$\hat{s}_n := \sum_{i=0}^n P_{\Omega}(\omega_i) \quad (22)$$

is necessarily convergent and thus a Cauchy sequence. This means that for any  $\hat{\epsilon} > 0$ , we can find an  $N_{\hat{\epsilon}} \in \mathbb{N}$  such that for all  $n > m > N_{\hat{\epsilon}}$  we have

$$\sum_{i=m}^n P_{\Omega}(\omega_i) \stackrel{P_{\Omega}(\omega_i) \geq 0}{=} |\hat{s}_n - \hat{s}_m| < \hat{\epsilon}. \quad (23)$$

Let  $\epsilon > 0$  be arbitrary but fixed. Choosing  $\hat{\epsilon} = \epsilon \frac{1-\gamma}{M}$  and  $n > m > N_{\hat{\epsilon}}$ , we indeed see that

$$|EG_{abs}^n - EG_{abs}^m| \stackrel{21}{\leq} \frac{M}{1-\gamma} \sum_{i=m}^n P_{\Omega}(\omega_i) \stackrel{23}{\leq} \frac{M}{1-\gamma} \hat{\epsilon} = \epsilon, \quad (24)$$

proving that  $(EG_{abs}^n)_n$  is indeed Cauchy and thus convergent in  $\mathbb{R}$ . This proves 19 for the one ordering of  $\Omega$  constructed, and thus 18 for *any* ordering of  $\omega \in \Omega$ . This proves the Lemma. □

In a similar vein, the expectations of the individual discounted rewards also exist for any finite MDP.

**Lemma 2.** *Let  $k \in \mathbb{N}$  be fixed but arbitrary. Then the expectation of the discounted reward at time step  $t+k$ , where  $t$  is the starting point in time for all trajectories considered, exist and are bounded by  $M$ . The same holds for the absolute value of the discounted reward at time step  $t+k$ . In other words, we have both  $\mathbb{E}_{\omega \sim P_\Omega}[\gamma^k r_{t+k}] \leq \gamma^k M$  and  $\mathbb{E}_{\omega \sim P_\Omega}[|\gamma^k r_{t+k}|] \leq \gamma^k M$ .*

*Proof.* The proof for either claim is nearly identical to arguments presented in the proof of Lemma 1. To avoid repetition we point out the two differences: instead of  $G_t(\omega)$ , we use either  $\gamma^k r_{t+k}$  or  $|\gamma^k r_{t+k}| = \gamma^k |r_{t+k}|$  depending on which of the two claims we want to prove. Secondly, the inequality in 21 is now justified by the assumption of universal boundedness 4, and the bound  $\frac{M}{1-\gamma}$  is to be replaced by  $\gamma^k M$ . Virtually every other argument can be copy pasted to prove both  $\mathbb{E}_{\omega \sim P_\Omega}[\gamma^k r_{t+k}] \leq \gamma^k M$  and  $\mathbb{E}_{\omega \sim P_\Omega}[|\gamma^k r_{t+k}|] \leq \gamma^k M$ .  $\square$

Before concluding this existential episode, let us present one more result and its application to our value function setting.

**Lemma 3.** *(Infinite linearity) Let  $P_\Omega$  be a distribution on the discrete trajectory space  $\Omega$ . Then*

$$\mathbb{E}_{\omega \sim P_\Omega} [G_t(\omega)] = \lim_{n \rightarrow \infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^n \gamma^k r_{t+k} \right] \quad (25)$$

*Proof.* The first step of the proof is to reduce the claim to that of linearity of the expectation for infinite random series, i.e. to reduce our claim to

$$\mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^{\infty} X_k \right] \stackrel{!}{=} \sum_{k=0}^{\infty} \mathbb{E}_{\omega \sim P_\Omega} [X_k] \quad (26)$$

and then refer to a standard argument in probability theory. To see that claim 25 can be rephrased in the form of 26, set  $X_k := \gamma^k r_{t+k}$  for  $k = 0, 1, \dots$ , giving equality of the respective left hand sides. Further, note that with these  $X_k$  and the fact that expectations are *always* linear for *finite* sums of random variables, we can see that the right hand side of 25 can be rewritten like

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^n \gamma^k r_{t+k} \right] &= \lim_{n \rightarrow \infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^n X_k \right] \\ &= \lim_{n \rightarrow \infty} \sum_{k=0}^n \mathbb{E}_{\omega \sim P_\Omega} [X_k] \\ &= \sum_{k=0}^{\infty} \mathbb{E}_{\omega \sim P_\Omega} [X_k], \end{aligned} \quad (27)$$

showing the claimed equivalence of both statements. To prove our modified claim 26, a standard probability theory argument states that all we require is for

$$\sum_{k=0}^{\infty} \mathbb{E}_{\omega \sim P_\Omega} [|X_k|] \stackrel{!}{<} \infty \quad (28)$$

to hold. This of course is equivalent to showing

$$\sum_{k=0}^{\infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \gamma^k |r_{t+k}| \right] \stackrel{!}{<} \infty, \quad (29)$$



where we already exploited  $\gamma > 0$ . We can quickly verify that indeed

$$\begin{aligned} \sum_{k=0}^n \mathbb{E}_{\omega \sim P_\Omega} [\gamma^k |r_{t+k}|] &\stackrel{Lem2}{\leq} \sum_{k=0}^n \gamma^k M \\ &< \sum_{k=0}^\infty \gamma^k M \\ &= \frac{M}{1-\gamma}. \end{aligned} \quad (30)$$

Letting  $n \rightarrow \infty$  on either side, we clearly have 29, thus proving the required condition and therefore proving the original claim 25.  $\square$

A useful consequence of this result is an alternative expression for our (action-) value functions  $V^\pi$  and  $Q^\pi$  for some policy  $\pi$  on a finite MDP.

**Corollary 3.** (*Limit your expectations*) *Let  $\pi$  be a policy on a finite MDP. Then*

1.

$$V^\pi(s) = \lim_{n \rightarrow \infty} \mathbb{E}_\pi \left[ \sum_{k=0}^n \gamma^k r_{t+k} | s_t = s \right] \quad (31)$$

2.

$$Q^\pi(s, a) = \lim_{n \rightarrow \infty} \mathbb{E}_\pi \left[ \sum_{k=0}^n \gamma^k r_{t+k} | s_t = s, a_t = a \right] \quad (32)$$

*Proof.* The proof for either identity goes all the way back (see the section surrounding equations 14, 15, 16 and 17) to the realisation that fixing the starting state  $s_t$  and then following a given policy  $\pi$  is inducing a distribution  $P_\Omega = P_{\Omega|s'}^\pi = P^\pi \circ P_{S,s'}$  on the trajectory space  $\Omega$ . We can write

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi [G_t | s_t = s] \\ &= \mathbb{E}_{\omega \sim P_\Omega} [G_t] \\ &\stackrel{Lem3}{=} \lim_{n \rightarrow \infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^n \gamma^k r_{t+k} \right] \\ &= \lim_{n \rightarrow \infty} \mathbb{E}_\pi \left[ \sum_{k=0}^n \gamma^k r_{t+k} | s_t = s \right], \end{aligned} \quad (33)$$

proving 31. Equivalently, fixing a starting state  $s_t$  and starting action  $a_t$  and then following a given policy  $\pi$  induces a distribution  $P_\Omega = P_{\Omega|s',a'}^\pi = P^\pi \circ P_{S,(s',a')}$ . This allows us to write

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}_\pi [G_t | s_t = s, a_t = a] \\ &= \mathbb{E}_{\omega \sim P_\Omega} [G_t] \\ &\stackrel{Lem3}{=} \lim_{n \rightarrow \infty} \mathbb{E}_{\omega \sim P_\Omega} \left[ \sum_{k=0}^n \gamma^k r_{t+k} \right] \\ &= \lim_{n \rightarrow \infty} \mathbb{E}_\pi \left[ \sum_{k=0}^n \gamma^k r_{t+k} | s_t = s, a_t = a \right], \end{aligned} \quad (34)$$

proving 32.  $\square$

### 3.2 Relationships

This subsection is dedicated to highlighting important and useful relationships between value and action-value functions. We will have a look at expressing one via the other, as well as recursive identities across time steps  $t$ .

Due to the heavy lifting done in the previous subsection as well as the inherently local (w.r.t time  $t$ , that is) behaviour of (finite) MDP, all probabilistic considerations in this section reduce to finite event spaces. In particular, instead of dealing with infinite but discrete trajectories  $\omega = (s_t, a_t, s_{t+1}, a_{t+1}, \dots)$ , we can focus on a finite set of state or action random variables  $(s_t, a_t), \dots, (s_{t+k}, a_{t+k})$ , where  $k$  will depend on the given context, but most of the times will not exceed 1. All expected values are therefore to be understood with respect to the finite distributions on these finite state-action space tuples resulting from following some specified policy  $\pi$ , sometimes with additional constraints on starting states or actions  $s_t, a_t$ . While the underlying trajectory distributions  $P_\Omega$  as outlined in the section surrounding equations ??, ??, ?? and ?? technically still apply, the properties defining a (finite) MDP as well as our previous work allow us to break free from this slightly technical "underground" world and instead work in the "above-ground" and easy-to-handle environment of finite state-action distributions.

The connection between these two worlds, in particular the transformation from trajectories  $\omega$  to individual states  $s_{t+k}$  and actions  $a_{t+k}$ , can be done by effectively grouping trajectories' probabilities that share, say, a state  $s$  at a fixed point in time  $t+k$ . To be precise, for some distribution  $P_\Omega$  on  $\Omega$ , we can generate a distributions  $P_{P_\Omega, S}$  for  $s_{t+k}$  on  $S$  by putting

$$P_{P_\Omega, S}(s_{t+k} = s) := \sum_{\omega \in \Omega_s^{t+k}} P_\Omega(\omega), \quad (35)$$

where  $\Omega_s^{t+k} := \{\omega = (s_t, a_t, s_{t+1}, a_{t+1}, \dots) \in \Omega \mid s_{t+k} = s\}$ . Similar constructions can be made to generate distributions over the action space  $A$  for the random variable  $a_{t+k}$  representing the action chosen at time  $t+k$ , and so on. In this sense, any policy  $\pi$  (with optional constraints on states and/or actions at some fixed point(s) in time) generates some  $P_\Omega$  on  $\Omega$ , which in turn generates a host of distributions on  $S$  and  $A$ , for individual states and actions, respectively. It is in this way that the following results can be connected to the results in the previous section. Personally, I found it helpful to keep in mind the concept of an "original source" distribution  $P_\Omega$  when dealing with different policies' behaviour and relative performance.

Having said all that, let us have a closer look at the recursive nature of both  $V^\pi$  and  $Q^\pi$  - a feature we will exploit throughout the rest of this transcript.

**Proposition 2.** (*Parametrized bellman equations*) *Let  $\pi$  be an arbitrary policy, and let  $V^\pi$  and  $Q^\pi$  its associated (action-)value functions. Then the following identities hold:*

1.  $V^\pi(s) = \sum_a \pi(s, a) \sum_{s'} P_{s, s'}^a (R_{s, s'}^a + \gamma V^\pi(s'))$
2.  $Q^\pi(s, a) = \sum_{s'} P_{s, s'}^a [R_{s, s'}^a + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a')]$

*Proof.* To prove the first identity consider

$$\begin{aligned}
V^\pi(s) &= \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s] \\
&= \mathbb{E}_\pi[r_t + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s] \\
&= \mathbb{E}_\pi[r_t | s_t = s] + \mathbb{E}_\pi[\gamma \sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s] \\
&\stackrel{2,3}{=} \sum_a \pi(s, a) \sum_{s'} P_{s,s'}^a R_{s,s'}^a + \gamma \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s] \\
&= \sum_a \pi(s, a) \sum_{s'} P_{s,s'}^a R_{s,s'}^a + \gamma \sum_a \pi(s, a) \sum_{s'} P_{s,s'}^a \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_{t+1} = s'] \\
&= \sum_a \pi(s, a) \sum_{s'} P_{s,s'}^a (R_{s,s'}^a + \gamma \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t'+k} | s_{t'} = s']) \\
&= \sum_a \pi(s, a) \sum_{s'} P_{s,s'}^a (R_{s,s'}^a + \gamma V^\pi(s')).
\end{aligned}$$

Note how we used the transitional probabilities and expected rewards to explicitly write out the expected value of the cumulative reward  $\sum_{k=0}^{\infty} \gamma^k r_{t+k}$  established in Lemma 1. Note also how we wrote  $\mathbb{E}_\pi$  to indicate that the expected value is w.r.t to a sequence of actions  $a_t, a_{t+1}, \dots$  and states  $s_t (= s, \text{ as per explicit condition } "s_t = s")$ ,  $s_{t+1}, \dots$  resulting from acting according to  $\pi$ , made explicit in subsequent steps including the term  $\pi(s, a)$ .

A similar line of thinking shows us that

$$\begin{aligned}
Q^\pi(s, a) &= \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s, a_t = a] \\
&= \mathbb{E}_\pi[r_t + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, a_t = a] \\
&= \mathbb{E}_\pi[r_t | s_t = s, a_t = a] + \gamma \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, a_t = a] \\
&\stackrel{2,3}{=} \sum_{s'} P_{s,s'}^a R_{s,s'}^a \\
&\quad + \gamma \sum_{s'} P_{s,s'}^a \sum_{a'} \pi(s, a') \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, a_t = a, s_{t+1} = s', a_{t+1} = a'] \\
&= \sum_{s'} P_{s,s'}^a R_{s,s'}^a + \gamma \sum_{s'} P_{s,s'}^a \sum_{a'} \pi(s, a') \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_{t+1} = s', a_{t+1} = a'] \\
&= \sum_{s'} P_{s,s'}^a R_{s,s'}^a + \gamma \sum_{s'} P_{s,s'}^a \sum_{a'} \pi(s, a') \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t'+k} | s_{t'} = s', a_{t'} = a'] \\
&= \sum_{s'} P_{s,s'}^a R_{s,s'}^a + \gamma \sum_{s'} P_{s,s'}^a \sum_{a'} \pi(s, a') Q^\pi(s', a') \\
&= \sum_{s'} P_{s,s'}^a (R_{s,s'}^a + \gamma \sum_{a'} \pi(s, a') Q^\pi(s', a'))
\end{aligned}$$

Note that we used the markov property which allowed us to drop past states and actions when going from line 4 to line 5. □

Remembering the definitions of  $P_{s,s'}^a$  and  $R_{s,s'}^a$ , these parametrizations can be rewritten in a slightly more compact way:

$$V^\pi(s) = \mathbb{E}_\pi[r_t + \gamma V^\pi(s_{t+1}) | s_t = s] \quad (36)$$

and

$$Q^\pi(s) = \mathbb{E}_\pi[r_t + \gamma Q^\pi(s_{t+1}, a_{t+1}) | s_t = s, a_t = a]. \quad (37)$$

These are the 'standard' bellman equations.

We now characterize the relationship between these two functions in the following

**Proposition 3.** (*QV Relationships*) *Let  $\pi$  be an arbitrary policy. Then the following identities hold:*

1.  $V^\pi(s) = \sum_a \pi(s, a) Q^\pi(s, a)$
2.  $V^\pi(s) = \sum_{a, s'} \pi(s, a) P_{s, s'}^a [R_{s, s'}^a + \pi(s', a') \gamma Q^\pi(s', a')]$
3.  $Q^\pi(s, a) = \sum_{s'} P_{s, s'}^a [R_{s, s'}^a + \gamma V^\pi(s')]$

*Proof.* To see that the first claims holds, we use the explicit distribution of taking an action  $a$  when in state  $s$  and following  $\pi$  to see that indeed

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s] \\ &= \sum_a \pi(s, a) \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s, a_t = a] \\ &= \sum_a \pi(s, a) Q^\pi(s, a). \end{aligned}$$

For the second claim, following a similar line of argument as we have done for Proposition 1, we see that

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi[r_t | s_t = s] + \gamma \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s] \\ &= \sum_a \sum_{s'} \pi(s, a) P_{s, s'}^a R_{s, s'}^a \\ &\quad + \sum_a \sum_{s'} \pi(s, a) P_{s, s'}^a \gamma \sum_{a'} \pi(s, a') \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, s_{t+1} = s', a_{t+1} = a'] \\ &= \sum_a \sum_{s'} \pi(s, a) P_{s, s'}^a R_{s, s'}^a \\ &\quad + \sum_a \sum_{s'} \pi(s, a) P_{s, s'}^a \gamma \sum_{a'} \pi(s, a') \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t'+k} | s_{t'} = s', a_{t'} = a'] \\ &= \sum_a \sum_{s'} \pi(s, a) P_{s, s'}^a (R_{s, s'}^a + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a')). \end{aligned}$$

The third equality uses the markov property. Lastly, we verify that

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}[r_t | s_t = s, a_t = a] + \gamma \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, a_t = a] \\ &= \sum_{s'} P_{s, s'}^a R_{s, s'}^a + \gamma \sum_{s'} P_{s, s'}^a \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+1+k} | s_t = s, a_t = a, s_{t+1} = s'] \\ &= \sum_{s'} P_{s, s'}^a R_{s, s'}^a + \gamma \sum_{s'} P_{s, s'}^a \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t'+k} | s_{t'} = s'] \\ &= \sum_{s'} P_{s, s'}^a (R_{s, s'}^a + \gamma V^\pi(s')) \end{aligned}$$

completing the proof. □

As before we give the more compact versions of these identities:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi[Q^\pi(s_t, a_t) | s_t = s], \\ V^\pi(s) &= \mathbb{E}_\pi[r_t + \gamma Q^\pi(s_{t+1}, a_{t+1}) | s_t = s] \end{aligned}$$

and

$$Q^\pi(s, a) = \mathbb{E}[r_t + \gamma V^\pi(s_{t+1}) | s_t = s, a_t = a].$$

Note that the last identity's expected value is *not* w.r.t  $\pi$  - that's because there is nothing random left to be determined according to  $\pi$ . The state  $s_t$  is given, the action  $a_t$

specified and fixed, and the expected value of the next state  $s_{t+1}$  is entirely dependent on how the environment reacts to this combination; and the term  $V^\pi$  is a deterministic function. This implies that policies sharing the same  $V$  also share the same  $Q$ . The reverse is not necessarily true. We formalize this realisation in the subsequent

**Corollary 4.** *Let  $\pi_1, \pi_2$  be two arbitrary policies such that  $V^{\pi_1} \equiv V^{\pi_2}$ . Then  $Q^{\pi_1} \equiv Q^{\pi_2}$*

*Proof.* This is most easily seen in the original, parametrized formulation of Proposition 2, 3. . Since both  $R_{s,s'}$  and  $P_{s,s'}^a$  are dependent on the environment only, and not on the policy in question, we clearly have

$$\begin{aligned} Q^{\pi_1}(s, a) &= \sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma V^{\pi_1}(s')] \\ &= \sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma V^{\pi_2}(s')] \\ &= Q^{\pi_2}(s, a) \end{aligned}$$

for all  $(s, a) \in S \times A$ . □

Another useful result derived from the same identity is formalized in the below

**Corollary 5.** *Let  $\pi$  be a policy for a finite MDP, and let  $(s, a) \sim P_{s,a}$  be a randomly distributed state-action pair. Then*

$$\mathbb{E}_{s,a \sim P_{s,a}}[Q^\pi(s, a)] = \mathbb{E}_{s_t, a_t \sim P_{s,a}}[r_t + \gamma V^\pi(s_{t+1})].$$

*Proof.* Since we are dealing with a finite MDP, both states and actions are drawn from a finite set  $S$  and  $A$ , respectively. We can therefore write

$$\begin{aligned} \mathbb{E}_{s,a \sim P_{s,a}}[Q^\pi(s, a)] &= \sum_{s,a} P_{s,a} Q^\pi(s, a) \\ &= \sum_{s,a} P_{s,a} \mathbb{E}[r_t + \gamma V^\pi(s_{t+1}) | s_t = s, a_t = a] \\ &= \mathbb{E}_{s_t, a_t \sim P_{s,a}}[r_t + \gamma V^\pi(s_{t+1})]. \end{aligned}$$

□

### 3.3 Comparing and improving

Now that we are a bit more comfortable with the concept of an (action-) value function, we can use it as a tool to quantify the *quality* of a given policy. Intuitively, it makes sense to regard a policy  $\pi_1$  that induces a higher expected reward when starting from a given state  $s$  than, say, another policy  $\pi_2$  as 'better' - at least for that given state. In other words, it makes sense to regard  $\pi_1$  as a better policy when starting from  $s$  than  $\pi_2$ , if and only if  $V^{\pi_1}(s) > V^{\pi_2}(s)$ . Expanding this intuitive measure of comparison beyond a single state  $s$  to *all* elements of  $S$ , we arrive at the following natural

**Definition 1.** (*Policy ranking*) *Let  $\pi_1, \pi_2$  be policies for a finite MDP. We say that*

$$\pi_1 \geq_V \pi_2$$

if and only if  $V^{\pi_1}(s) \geq V^{\pi_2}(s)$  for all  $s \in S$ . We say that

$$\pi_1 >_V \pi_2$$

if and only if  $\pi_1 \geq_V \pi_2$  and there exists at least one  $s \in S$  such that  $V^{\pi_1}(s) > V^{\pi_2}(s)$ . Finally, we say that

$$\pi_1 =_V \pi_2$$

if and only if  $V^{\pi_u}(s) = V^{\pi_l}(s)$  for all  $s \in S$ .

This ranking of policies w.r.t  $\geq_V$  induces a partial ordering on the set of policies  $\Pi$ . Note that it is possible that neither  $\pi_1 \geq_V \pi_2$  nor  $\pi_1 \leq_V \pi_2$  for a given pair of policies  $\pi_1, \pi_2$ , since we demand that one value function exceeds the other for *all*  $s \in S$ . In other words,  $\geq_V$  really only is a *partial* ordering on the set of policies  $\Pi$ .

We sometimes informally refer to the relation  $\pi_1 \geq_V \pi_2$  as ' $\pi$  is *at least as good as*  $\pi_2$ ', to  $\pi_1 >_V \pi_2$  as ' $\pi_1$  is *better than / an actual improvement over*  $\pi_2$ ' and to  $\pi_1 =_V \pi_2$  as ' $\pi_1$  is *as good as / equally good as*  $\pi_2$ '.

We have, even at this early stage, established enough theory to characterize some cases where a direct comparison of policies w.r.t  $\geq_V$  is possible.

**Theorem 1.** (*Policy improvement theorem*) Let  $\pi_u, \pi_l$  be two different policies for a finite MDP such that

$$\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] \geq V^{\pi_l}(s)$$

for all  $s \in S$ . Then

$$\pi_u \geq_V \pi_l.$$

Before we begin the proof, let us formulate the above statement in a slightly less formal way. Our condition on  $\pi_u$  and  $\pi_l$  can be paraphrased as follows: If the expected reward generated by following  $\pi_u$  for *one* time step (note the expected value is indexed with  $\pi_u$ , indicating that the one remaining free random variable  $a_t$  is chosen according to  $\pi_u$  *conditioned*, i.e. fixed in its state variable, on the value of the state  $s$ ) and then following  $\pi_l$  for all subsequent time steps is *always* (i.e. for every starting state  $s$ ) greater or equal to the expected reward generated by following  $\pi_l$  *from the start*, then the policy  $\pi_u$  must be at least as good as  $\pi_l$  overall. In other words, if 'prepending' your actions with one action from a specified policy does not deteriorate rewards, the policy generating that one inserted action at the start of your journeys is at least as good as the policy being prepended. We will actually use this idea in an induction approach to show that, as we iteratively increase the number of time steps in which the actions are being chosen according to  $\pi_u$  before switching back to  $\pi_l$ , the expected reward does not decrease as well as converges to  $V^{\pi_u}$ .

Another thing to note is that, if the policy  $\pi_u$  is deterministic, our condition in the theorem reduces to

$$Q^{\pi_l}(s, \pi_u(s)) \geq V^{\pi_l}(s)$$

as the expected value of a constant random variable reduces to that constant value.

*Proof.* We first need to extend our assumption to the case where the state  $s$  appearing on both sides is not fixed, but more generally a random variable distributed according to, say, some distribution  $P_s$ . Since  $P_s$  is a distribution over finite states, we can see that indeed

$$\begin{aligned}\mathbb{E}_{s \sim P_s}[V^{\pi_l}(s)] &= \sum_s P_s V^{\pi_l}(s) \\ &\leq \sum_s P_s \mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] \\ &= \sum_s P_s \sum_a \pi_u(s, a) Q^{\pi_l}(s, a) \\ &= \mathbb{E}_{s_t, a_t \sim (P_{s_t}, \pi_u(s_t, \cdot))}[Q^{\pi_l}(s_t, a_t)].\end{aligned}$$

The second and main part of the proof will consist of showing the aforementioned policy improvement via prepending  $\pi_l$  with  $\pi_u$ . Formally, this means that for any  $s \in S$ ,  $n = 1, 2, \dots$  the inequality

$$V^{\pi_l}(s) \stackrel{!}{\leq} \mathbb{E}_{\pi_u} \left[ \sum_{k=0}^{n-1} \gamma^k r_{t+k} | s_t = s \right] + \gamma^n \mathbb{E}_{\pi_u} [V^{\pi_l}(s_{t+n}) | s_t = s]$$

holds. We will show this claim by induction over  $n$ .

For the induction start, let  $s \in S$  be arbitrary but fixed. We then see that, by our assumption, the definition of the value function  $V^\pi$ , and Proposition 2, 3., we have

$$\begin{aligned}V^{\pi_l}(s) &\leq \mathbb{E}_{\pi_u}[Q^{\pi_l}(s_t, a_t) | s_t = s] \\ &= \sum_a \pi_u(s, a) Q^{\pi_l}(s, a) \\ &= \sum_a \pi_u(s, a) \mathbb{E}[r_t + \gamma V^{\pi_l}(s_{t+1}) | s_t = s, a_t = a] \\ &= \mathbb{E}_{\pi_u}[r_t + \gamma V^{\pi_l}(s_{t+1}) | s_t = s] \\ &= \mathbb{E}_{\pi_u}[r_t | s_t = s] + \gamma \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+1}) | s_t = s],\end{aligned}$$

proving the claim for  $n = 1$ . Remember that the index  $\pi_u$  denotes that any implicit intermediate action  $a$  was taken according to  $\pi_u$ .

For the induction step, let us introduce some additional notation. Let  $s_{t+k} \sim P_s^{k * \pi_u}$  denote the distribution for the state at  $t+k$  given that the state at time  $t$  was  $s$  and the subsequent  $k$  action(s)  $a_t, \dots, a_{t+k-1}$  were chosen according to  $\pi_u$ . Then we can apply our expected value version of the initial assumption to see that

$$\begin{aligned}V^{\pi_l}(s) &\stackrel{I.S.}{\leq} \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] + \gamma^{n-1} \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+n-1}) | s_t = s] \\ &= \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] + \gamma^{n-1} \mathbb{E}_{s_{t+n-1} \sim P_s^{(n-1) * \pi_u}}[V^{\pi_l}(s_{t+n-1})] \\ &\leq \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] \\ &\quad + \gamma^{n-1} \mathbb{E}_{(s_{t+n-1}, a_{t+n-1}) \sim (P_s^{(n-1) * \pi_u}, \pi_u(s_{t+n-1}, \cdot))}[Q^{\pi_l}(s_{t+n-1}, a_{t+n-1})]\end{aligned}$$

Applying Corollary 2 with  $P_{s', a'} = (P_s^{1 * \pi_u}, \pi_u(s', \cdot))$  to the last expression in the above sequence, we arrive at

$$\begin{aligned}
V^{\pi_l}(s) &\leq \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] \\
&\quad + \gamma^{n-1} \mathbb{E}_{(s_{t+n-1}, a_{t+n-1}) \sim (P_s^{(n-1)*\pi_u}, \pi_u(s_{t+n-1}, \cdot))} [r_{t+n-1} + \gamma V^{\pi_l}(s_{t+n})] \\
&= \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] \\
&\quad + \gamma^{n-1} \mathbb{E}_{\pi_u}[r_{t+n-1} + \gamma V^{\pi_l}(s_{t+n}) | s_t = s] \\
&= \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-2} \gamma^k r_{t+k} | s_t = s] + \mathbb{E}_{\pi_u}[\gamma^{n-1} r_{t+n-1} | s_t = s] \\
&\quad + \gamma^n \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+n}) | s_t = s] \\
&= \mathbb{E}_{\pi_u}[\sum_{k=0}^{n-1} \gamma^k r_{t+k} | s_t = s] + \gamma^n \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+n}) | s_t = s],
\end{aligned}$$

completing the induction step and proving the claim. Letting  $n \rightarrow \infty$ , we see that

$$\begin{aligned}
V^{\pi_l}(s) &\leq \mathbb{E}_{\pi_u}[\sum_{k=0}^n \gamma^k r_{t+k} | s_t = s] + \gamma^n \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+n}) | s_t = s] \\
&\xrightarrow{n \rightarrow \infty} V^{\pi_u}(s) + 0.
\end{aligned}$$

For the second term, we have implicitly used that

$$0 \stackrel{\infty \leftarrow n}{\leftarrow} \min_{s \in S} V^{\pi_l}(s) \leq \gamma^{n_a} \mathbb{E}_{\pi_u}[V^{\pi_l}(s_{t+n_a}) | s_t = s] \leq \gamma^{n_a} \max_{s \in S} V^{\pi_l}(s) \xrightarrow{n \rightarrow \infty} 0.$$

The convergence claim regarding the first term was proven in Corollary 3, 31.

Informally speaking, taking an infinite number of actions according to  $\pi_u$  before switching to  $\pi_l$  essentially means simply following  $\pi_u$  and generates, independent of the starting state  $s$ , an expected cumulative reward that is at least as high as the one generated by simply following  $\pi_l$ . This proves the theorem.  $\square$

The above theorem states that  $\pi_u \geq_V \pi_l$  holds whenever  $\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] \geq V^{\pi_l}(s)$  for all  $s \in S$ , then  $\pi_u \geq_V \pi_l$ . Can we tweak these assumptions to guarantee actual improvement over  $\pi_l$ , i.e.  $\pi_u >_V \pi_l$ ?

**Theorem 2.** (*Improved policy improvement theorem*) Let  $\pi_u, \pi_l$  be two different policies for a finite MDP.

1. If

$$\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] = V^{\pi_l}(s)$$

for all  $s \in S$ , then

$$\pi_u =_V \pi_l.$$

2. If

$$\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] \geq V^{\pi_l}(s)$$

for all  $s \in S$ , then

$$\pi_u \geq_V \pi_l.$$



3. If

$$\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] \geq V^{\pi_l}(s)$$

for all  $s \in S$ , then

$$\pi_u \succ_V \pi_l$$

and there is at least one  $s \in S$  such that

$$\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] > V^{\pi_l}(s),$$

then

$$\pi_u \succ_V \pi_l.$$

*Proof.* Claim 2. is just the original Policy Improvement Theorem.

To see claim 1., simply replace the  $\geq$  in the proof of the Policy Improvement Theorem with  $=$ . This equality persists through the induction step for all starting states  $s$  and yields equality of the respective value functions  $V^{\pi_u}$  and  $V^{\pi_l}$ .

To see the third claim, we again use the Policy Improvement Theorem to see that  $\pi_u \succ_V \pi_l$ . We now have to find at least one  $s \in S$  satisfying  $V^{\pi_u}(s) > V^{\pi_l}(s)$ . Repeating the argument displayed in the proof of the Policy Comparison Theorem, we see that, applying it to the one state  $s$  guaranteed by our starting assumption to achieve  $\mathbb{E}_{a \sim \pi_u(s, \cdot)}[Q^{\pi_l}(s, a)] > V^{\pi_l}(s)$ , the strict inequality persists throughout the induction step and really yields  $V^{\pi_u}(s) > V^{\pi_l}(s)$ . This shows  $\pi_l u \succ_V \pi_l$  and completes the proof.  $\square$

We have now established some degree of comparability based on the policies' in question's  $V$  and  $Q$  functions. Before we can use this result to iteratively construct better and better policies, let us formalize the concept of a greedy policy.

**Definition 2.** Let  $\pi$  be a policy for a finite MDP and let  $Q$  be some (possibly but not necessarily  $\pi$ 's) action-value function. For any state  $s \in S$ , denote by

$$\text{maxact}^Q(s) := \{a \in A \mid Q(s, a) = \max Q(s, a)\},$$

the set of actions at which  $Q(s, \cdot)$  achieves its maximum. We say that  $\pi$  is  $Q$ -greedy if and only if for every  $s \in S$  we have

$$\sum_{a \in A} \pi(s, a) = \sum_{a \in \text{maxact}^Q(s)} \pi(s, a) = 1$$

If a policy is greedy w.r.t its own  $Q$  function, we simply call it *greedy*. In other words, a *greedy* policy  $\pi$  only chooses among the actions that maximise its own  $Q^\pi$ -function for the given state. The only thing that matters for these policies when making a decision in state  $s$  on what to do next is whether that immediate action achieves a maximal 'immediate' pay-off of  $\max_{a \in A} Q^\pi(s, a)$ .

At first glance, such a policy might be somewhat short-sighted, seemingly ignoring potential future consequences of its immediate actions for the sake of instantaneous profit.

However, the notion of the action-value function is to encode the future cumulative reward's expected value - in other words, it is a function that very much 'looks ahead' and considers future consequences; namely, all of them.

This means that a *greedy* policy might not be as shortsighted, and therefore not such a bad thing, after all. Indeed, we will later see that the best policies are exactly the ones that are *greedy*.

Let us therefore have a closer look at a greedy policy's value function.

**Corollary 6.** (*Greedy policy values*) Let  $\pi_g$  be a greedy policy for a finite MDP. Then

$$V^{\pi_g} \equiv \max_{a \in A} Q^{\pi_g}(\cdot, a).$$

*Proof.* We write for any fixed but arbitrary  $s \in S$

$$\begin{aligned} V^{\pi_g}(s) &= \sum_a \pi_g(s, a) Q^{\pi_g}(s, a) \\ &= \sum_{a \in \text{maxact}^{Q^{\pi_g}}(s)} \pi_g(s, a) Q^{\pi_g}(s, a) \\ &= \max_{a \in A} Q^{\pi_g}(s, a). \end{aligned}$$

□

In other words, a greedy policy's value function is just the maximum of its action-value function at the present state, taken over all possible actions. That is not surprising, seeing as maximising its own  $Q$  function is what is driving a greedy policy as per definition.

The following result shows that any group of policies greedy w.r.t some  $Q$  function are of the same quality. Unsurprisingly, the identical strategy of local maximisation does not lead to great deal of variability in policies' performances. It also shows that any policy that is not *greedy* can be improved by making it *greedy*.

**Lemma 4.** (*Greedy policy improvement*) Let  $\pi_g$  and  $\pi_c$  be policies for finite MDP, and let  $\pi_g$  be  $Q^{\pi_c}$ -greedy. If  $\pi_c$  is  $Q^{\pi_c}$ -greedy, too, then  $\pi_g =_V \pi_c$ . Otherwise,  $\pi_g >_V \pi_c$ .

*Proof.* Let  $s \in S$  be fixed but arbitrary. We write

$$\begin{aligned} \mathbb{E}_{a \sim \pi_g(s, \cdot)}[Q^{\pi_c}(s, a)] &= \sum_a \pi_g(s, a) Q^{\pi_c}(s, a) \\ &= \sum_{a \in \text{maxact}^{\pi_g}(s)} \pi_g(s, a) \max_{a \in A} Q^{\pi_c}(s, a) \\ &= \max_{a \in A} Q^{\pi_c}(s, a) \\ &\geq \sum_{a \in A} \pi_c(s, a) Q^{\pi_c}(s, a) \\ &= V^{\pi_c}(s). \end{aligned}$$

If  $\pi_c$  is  $Q^{\pi_c}$ -greedy, then we have equality in the above chain for all  $s \in S$  and the Improved Policy Improvement Theorem, 1., implies  $\pi_g =_V \pi_c$ . If  $\pi_c$  is not  $Q^{\pi_c}$ -greedy, then there is at least one  $s \in S$  such that  $>$  holds (and still  $\geq$  for all other  $s$ ). In that case the Improved Policy Improvement Theorem, 3., yields the desired claim. □

Inspecting the above proof we notice an important detail: Any deterministic policy that chooses an action from  $\text{maxact}^{\pi_c}(s)$  given any state  $s$  satisfies the condition of the corollary and thus is at least as good as (the potentially probabilistic)  $\pi_c$ . This means that for every probabilistic policy there is a deterministic policy that is at least as good, and it can be constructed explicitly by letting it choose any one action that maximises the probabilistic policy's action-value function.

## 4 Optimal policies and how to find them

In this section we will use our concept of policy comparison to formalize and analyse the concept of a *best* policy. We will investigate existence, uniqueness and characterizations of such policies. We will mainly build on the previous sections results, but will at a later point be forced to introduce some useful tools from functional analysis.

It is fair to say that most, if not all, of the subsequent results presented in this script deal with the analysis of optimal policies. In particular, we will formally answer the following questions:

- Is there always a (unique) optimal policy for a finite MDP?
- Are there any characteristic traits that all optimal policies share, and if so, how can we make use of them to find these policies?
- Can we give a constructive way of finding or at least approximating these optimal policies?

With this road map in mind, let us start our journey by clarifying what we mean by an *optimal* or *best* policy.

### 4.1 Existential crisis

Let us begin this subsection by formalizing the concept of an optimal policy.

**Definition 3.** (*Optimal policy*) Let  $\pi^*$  be a policy for a finite MDP. We call  $\pi^*$  an *optimal policy* if and only if we have

$$\pi^* \geq_V \pi$$

for all policies  $\pi \in \Pi$ .

In other words, a policy  $\pi^*$  whose value function  $V^{\pi^*}$  dominates the value functions  $V^\pi$  of all other policies is optimal for the given finite MDP, and what we consider 'best'. This approach seems sensible, since optimality of a policy according to the above definition maximises the expected cumulative reward (that is what a policy's value function represents).

It is not at all clear, or even intuitive, that such an optimal policy need exist for a given finite MDP. As we will see later in this section, an optimal policy indeed does exist under the conditions treated in this article. Somewhat unintuitively, will first analyse the behaviour of surmised optimal policies, and use these properties to show to then prove their existence. This approach brings us to the next subsection.

## 4.2 Behavioural issues

This subsection is dedicated to deriving some useful characteristics of the mysterious (as of yet) optimal policy of a given finite MDP. These results will turn out to be essential for the construction of explicit algorithms aiming to converge on optimal policies. We start of with an easy corollary that states that, while there might be more than one optimal policy, they all share the same value and action value function. Since it is the value function that we use to rank and disambiguate policies from one another, the following statement essentially says that there is no need to tell one optimal policy apart from any other optimal policy, they are - literally for what it's worth - the same thing.

**Corollary 7.** *(It's all the same thing) Let  $\pi_1, \pi_2$  be two optimal policies for a finite MDP. Then*

1. 
$$V^{\pi_1} \equiv V^{\pi_2} \tag{38}$$

2. 
$$Q^{\pi_1} \equiv Q^{\pi_2} \tag{39}$$

*Proof.* Since both  $\pi_1$  and  $\pi_2$  are optimal policies, we have both

$$\pi_1 \geq_V \pi_2 \Leftrightarrow V^{\pi_1}(s) \geq V^{\pi_2}(s) \quad \forall s \in S \tag{40}$$

as well as

$$\pi_2 \leq_V \pi_1 \Leftrightarrow V^{\pi_2}(s) \geq V^{\pi_1}(s) \quad \forall s \in S, \tag{41}$$

clearly implying the equality stated in 38. Claim 39 follows from 38 and Corollary 4.  $\square$

Knowing that all optimal policies share the same value and action value function, we can define the optimal versions of these functions independently of an attached optimal policy  $\pi^*$ .

**Definition 4.** *(Optimal (action-)values) Let  $\pi^*$  be any optimal policy for a given finite MDP. We call  $V^* = V^{\pi^*}$  the optimal value function, and  $Q^* = Q^{\pi^*}$  the optimal action-value function.*

Corollary 7 guarantees the welldefinedness of the above terms. We go on to show a useful property of the optimal value function  $V^*$ .

**Lemma 5.** *(Different kind of optimum) For each  $s \in S, a \in A$ , the following hold:*

1. 
$$V^*(s) = \max_{\pi \in \Pi} V^\pi(s) \tag{42}$$

2.

$$Q^*(s, a) = \max_{\pi \in \Pi} Q^\pi(s, a) \quad (43)$$

In other words, both the optimal value function and the optimal action-value function can be seen as the results of pointwise optimization over the state space  $S$ .

*Proof.* We start with the first claim. Let  $s \in S$  be fixed but arbitrary. By the very definition of  $V^*$  as the value function of an optimal policy  $\pi^*$  which satisfies ??, we have

$$\max_{\pi \in \Pi} V^\pi(s) = \max_{\pi^* \in \Pi, \pi^{\text{optimal}}} V^{\pi^*}(s). \quad (44)$$

By Corollary 7 we can *arbitrarily* choose any optimal policy  $\pi_1^*$  to obtain

$$\max_{\pi^* \in \Pi, \pi^{\text{optimal}}} V^{\pi^*}(s) = V^{\pi_1^*}(s) = V^*(s), \quad (45)$$

proving the identity 42. As for the second claim, let  $(s, a) \in S \times A$  be arbitrary but fixed. Then

$$\begin{aligned} \max_{\pi \in \Pi} Q^\pi(s, a) &\stackrel{\text{Prop42}}{=} \max_{\pi \in \Pi} \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')] \\ &= \sum_{s'} P_{ss'}^a R_{ss'}^a + \max_{\pi \in \Pi} \sum_{s'} P_{ss'}^a \gamma V^\pi(s'). \end{aligned} \quad (46)$$

where we have used that the term quantifying the expectation of the immediate reward  $\mathbb{E}[r_t | s_t = s, a_t = a] = \sum_{s'} P_{ss'}^a R_{ss'}^a$  is independent of the particular policy  $\pi$  if both  $s_t$  and  $a_t$  are fixed at  $s$  and  $a$ , respectively. We see that

$$\begin{aligned} &\sum_{s'} P_{ss'}^a R_{ss'}^a + \max_{\pi \in \Pi} \sum_{s'} P_{ss'}^a \gamma V^\pi(s') \\ &\leq \sum_{s'} P_{ss'}^a R_{ss'}^a + \sum_{s'} P_{ss'}^a \gamma \max_{\pi \in \Pi} V^\pi(s') \\ &\stackrel{42}{=} \sum_{s'} P_{ss'}^a R_{ss'}^a + \sum_{s'} P_{ss'}^a \gamma V^*(s') \\ &= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')]. \end{aligned} \quad (47)$$

Arbitrarily choosing an optimal policy  $\pi^*$  we finally see that

$$\begin{aligned} &\sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')] \\ &\stackrel{\text{Def4}}{=} \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^{\pi^*}(s')] \\ &\stackrel{\text{Prop3}}{=} Q^{\pi^*}(s, a) \\ &\stackrel{\text{Def4}}{=} Q^*(s, a), \end{aligned} \quad (48)$$

proving " $\geq$ " in 43. To see that " $\leq$ " holds, consider, again using some arbitrary optimal policy  $\pi^*$ , that, since the optimal policy  $\pi^*$ , too, is an element of the set of *all* policies  $\Pi$ , we also have

$$Q^*(s, a) \stackrel{\text{Def4}}{=} Q^{\pi^*}(s, a) \leq \max_{\pi \in \Pi} Q^\pi(s, a). \quad (49)$$

Since  $(s, a) \in S \times A$  was arbitrary,  $\geq$  and  $\leq$  together yield 43.  $\square$

It is worth reminding ourselves that this property assumes the existence of an optimal policy  $\pi^*$  with corresponding value and action value functions  $V^{\pi^*} = V^*$ ,  $Q^{\pi^*} = Q^*$ , respectively. In other words, the right hand sides of both identities are well defined, but there is no guarantee that the policy counterpart, as implied by the left hand sides, exists.

We next show that greedyness is a necessary condition for a policy to be optimal.

**Lemma 6.** (*Necessary conditions for optimal policies*) *Let  $\pi^*$  be an optimal policy. Then  $\pi^*$  is greedy.*

*Proof.* Assume that  $\pi^*$  isn't greedy. By setting

$$\pi_{imp}(s, a) := \begin{cases} \frac{1}{|\maxact^{Q^{\pi^*}}(s)|}, & a \in \maxact^{Q^{\pi^*}} \\ 0, & a \notin \maxact^{Q^{\pi^*}}, \end{cases} \quad (50)$$

we can see that  $\pi_{imp}$  and  $\pi^*$  satisfy the conditions of the Improved Policy improvement Theorem 2, implying that  $\pi_{imp}$  is  $Q^{\pi^*}$  greedy. Applying Lemma 4 therefore yields

$$\pi_{imp} >_V \pi^*, \quad (51)$$

clearly a violation of our optimality assumption on  $\pi^*$  requiring

$$V^{\pi^*}(s) \geq V^{\pi}(s) \quad (52)$$

for all  $s \in S$ ,  $\pi \in \Pi$ . Hence,  $\pi^*$  must be greedy.  $\square$

For some  $\pi \in \Pi$ , consider the original *Bellman optimality equation* given by

$$Q^{\pi}(s, a) \stackrel{!}{=} \mathbb{E}[r_t + \gamma \max_{a'} Q(s_{t+1}, a') | s_t = s, a_t = a]. \quad (53)$$

We next show that a policy  $\pi$  satisfies this identity holds true for greedy policies.

**Lemma 7.** (*Bell optimality equations*) *Let  $\pi \in \Pi$  be an arbitrary policy. If  $\pi$  is greedy then it satisfies the Bellman optimality equation 53.*

*Proof.* Let  $\pi$  be greedy. By one of our early results, namely Proposition 3 and Corollary 6, we clearly have

$$\begin{aligned} Q^{\pi}(s, a) &= \mathbb{E}[r_t + \gamma V^{\pi}(s_{t+1}) | s_t = s, a_t = a] \\ &= \mathbb{E}[r_t + \gamma \max_{a' \in A} Q^{\pi}(s_{t+1}, a') | s_t = s, a_t = a]. \end{aligned} \quad (54)$$

$\square$

### 4.3 The Bellman operator

In this section, we will develop a slightly different view point on the action-value function of a policy - in particular, the action value function of greedy policies. We ask the reader to bear with us on this seemingly sudden deviation from our path in the hope that these results' benefits will become apparent soon enough. To this end, consider the set of non-negative valued functions mapping into  $\mathbb{R}^+$ , that is, functions of the form

$$\mathbb{Q}_{S,A} := \{Q : (S, A) \mapsto \mathbb{R}_0^+\} \quad (55)$$

equipped with the maximum norm

$$|Q|_\infty := \max_{(s,a) \in S \times A} |Q(s, a)|. \quad (56)$$

Since both state and action states are finite, respectively, we can choose an arbitrary ordering of actions and states  $A = \{a_1, \dots, a_{|A|}\}$  and  $S = \{s_1, \dots, s_{|S|}\}$ . This in turn allows us to put an ordering on the state-action pairs  $(s, a)_1 = (s_1, a_1), (s, a)_2 = (s_1, a_2), \dots, (s, a)_{n_{S,A}} = (s_{|S|}, a_{|A|})$ , where  $n_{S,A} := |S| \cdot |A|$ . Using this ordering, we can embed our set of action-value functions  $\mathbb{Q}_{S,A}$  into the normed space  $(\mathbb{R}^{n_{S,A}}, |\cdot|_\infty)$  via the canonical mapping

$$\begin{aligned} E : \mathbb{Q}_{S,A} &\hookrightarrow \mathbb{R}^{n_{S,A}} \\ Q &\mapsto (Q((s, a)_1), \dots, Q((s, a)_{n_{S,A}})). \end{aligned} \quad (57)$$

Thus, we can interpret our action-value functions as a subset of the Banach space  $(\mathbb{R}^{n_{S,A}}, |\cdot|_\infty)$ . It can easily be shown that  $\mathbb{Q}_{S,A} \subset \mathbb{R}^{n_{S,A}}$  is a closed subset w.r.t  $|\cdot|_\infty$ . Since closed subspaces of complete spaces are complete themselves w.r.t the parent norm, we have obtained the following

**Lemma 8.** *Let  $\mathbb{Q}_{S,A}$  be defined as in equation 55. Then  $(\mathbb{Q}_{S,A}, |\cdot|_\infty)$  is a Banach space.*

The above result provides the perfect frame work to revisit Bellman's optimality equation from an operator theory point of view. Consider the non-linear operator on  $\mathbb{Q}_{S,A}$  defined by

$$\begin{aligned} B_{\mathbb{Q}} : \mathbb{Q}_{S,A} &\rightarrow \mathbb{Q}_{S,A} \\ Q &\mapsto B_{\mathbb{Q}} \end{aligned} \quad (58)$$

where the image  $B_{\mathbb{Q}}(Q)$  of  $Q$  under  $B_{\mathbb{Q}}$  is given by

$$B_{\mathbb{Q}}(Q)(s, a) := \mathbb{E}[r_t + \gamma \max_{a'} Q(s_{t+1}, a') | s_t = s, a_t = a], \quad (59)$$

$s \in S, a \in A$ . Note that  $B_{\mathbb{Q}}$  is well-defined due to  $r_t \geq 0$  for any sensible  $t$ . It is also worth noting that for  $B_{\mathbb{Q}}$  to be applied, no policy  $\pi$  needs to be specified. Indeed,  $Q$  does not even need to be an actual action-value function derived from some policy  $\pi$  (although it usually will be - see below sections). All it requires is a function mapping state-action pairs to the non-negative part of the real numbers.



What does any of this have to do with our quest of finding and constructing (and, not to forget, proving the existence of) some optimal policy  $\pi^*$ ? Well, we know that greedyness is a necessary condition for optimality (Lemma 6), and that satisfying the Bellman optimality equation is a necessary condition for greedyness (Lemma 7). To extend this reasoning by one more step, we can see that satisfying the Bellman optimality equation ?? is, by construction of  $B_{\mathbb{Q}}$ , equivalent to  $Q^\pi = B_{\mathbb{Q}}(Q^\pi)$  in  $\mathbb{Q}_{S,A}$ . Thus we see the following implications hold true:

$$\pi^* \text{ optimal} \implies \pi^* \text{ greedy} \implies Q^{\pi^*} \text{ is a fixed point of } B_{\mathbb{Q}}. \quad (60)$$

A direct consequence of this is: if there is no fixed point of  $B_{\mathbb{Q}}$ , there can't be any optimal policy  $\pi^*$ . The good news is that not only does such a fixed point always exist, but it is also unique.

**Proposition 4.** *Let  $S$  and  $A$  be finite state and action spaces, respectively. Then the operator  $B_{\mathbb{Q}}$  has a unique fixed point  $Q^*$ .*

A word of warning: Our notation  $Q^*$  heavily implies that we are dealing with the optimal policy's action value function. *We have not shown this yet.* The implications in 60 are strictly one-way, and there is still some work to be done to link the fixed point of  $B_{\mathbb{Q}}$  back to some optimal policy  $\pi^*$ .

*Proof.* We will prove the claim by showing that  $B_{\mathbb{Q}}$  is a contraction on  $(\mathbb{Q}_{S,A}, |\cdot|_\infty)$  and then apply the Banach Fixed Point Theorem (see appendix). In other words, we need to show

$$|B_{\mathbb{Q}}(Q_1) - B_{\mathbb{Q}}(Q_2)|_\infty < c |Q_1 - Q_2|_\infty \quad (61)$$

for some  $0 < c < 1$  and any pair  $Q_1, Q_2 \in \mathbb{Q}_{S,A}$ . Noting that we can rewrite

$$\begin{aligned} B_{\mathbb{Q}}(Q)(s, a) &= \mathbb{E}[r_t + \gamma \max_{a'} Q(s_{t+1}, a') | s_t = s, a_t = a] \\ &= \sum_{s' \in S} P_{s,s'}^a (R_{s,s'}^a + \gamma \max_{a' \in A} Q(s', a')), \end{aligned} \quad (62)$$

we see that

$$\begin{aligned}
|B_{\mathbb{Q}}(Q_1) - B_{\mathbb{Q}}(Q_2)|_{\infty} &= \max_{(s,a) \in S \times A} \left| \sum_{s' \in S} P_{s,s'}^a \gamma \left( \max_{a' \in A} Q_1(s', a') - \max_{a' \in A} Q_2(s', a') \right) \right| \\
&\leq \gamma \max_{(s,a) \in S \times A} \left( \sum_{s' \in S} P_{s,s'}^a \left| \max_{a' \in A} Q_1(s', a') - \max_{a' \in A} Q_2(s', a') \right| \right) \\
&\leq \gamma \max_{(s,a) \in S \times A} \left( \sum_{s' \in S} P_{s,s'}^a \max_{a' \in A} |Q_1(s', a') - Q_2(s', a')| \right) \\
&\leq \gamma \max_{(s,a) \in S \times A} \left( \sum_{s' \in S} P_{s,s'}^a \max_{(s'', a') \in S \times A} |Q_1(s'', a') - Q_2(s'', a')| \right) \\
&= \gamma \max_{(s,a) \in S \times A} \left( \sum_{s' \in S} P_{s,s'}^a |Q_1 - Q_2|_{\infty} \right) \\
&= \gamma |Q_1 - Q_2|_{\infty}.
\end{aligned} \tag{63}$$

Since  $0 < \gamma < 1$ , there is a  $\gamma < c < 1$  such that equation 61 holds, proving  $B_{\mathbb{Q}}$  is a contraction on  $(\mathbb{Q}_{S,A}, |\cdot|_{\infty})$ . By the Banach Fixed Point Theorem it has a unique fixed point  $Q^* \in \mathbb{Q}_{S,A}$ .  $\square$

#### 4.4 Building optimal policies

We are finally able to formulate this article's main result.

**Theorem 3.** (*Building optimal policies*) Let  $\pi_0 \in \Pi$  be an arbitrary (deterministic or probabilistic) policy for a given finite MDP. For  $i = 1, \dots$ , iteratively put

$$\pi_i(s, a) := \begin{cases} 1, & a = a_i, i = \min\{j | a_j \in \maxact^{Q^{\pi_{i-1}}}(s)\} \\ 0, & \text{otherwise.} \end{cases} \tag{64}$$

Then there exists an  $N \in \mathbb{N}$  such that the following hold:

1.  $\pi_{i-1} <_V \pi_i$  for  $i = 1, \dots, N$
2.  $\pi_N =_V \pi_j$  for  $j = N, \dots$ ,
3.  $\pi_N = \pi^*$ , i.e.  $\pi_N$  is optimal

The notation of the first case in 64 refers to the ordering of the states  $(a_1, \dots, a_{|A|}) = A$  already mentioned earlier. While the ordering is arbitrary and not required, we refer to it here to formulate our policy improvement algorithm in a closed, succinct form. What we mean to say is to really only chose the 'first' (w.r.t our arbitrary ordering on  $A$ ) action in  $\maxact^{Q^{\pi_{i-1}}}(s)$  and to ignore all other actions, even if  $|\maxact^{Q^{\pi_{i-1}}}(s)| > 1$ . This way, we remain in the space of *deterministic* policies.

*Proof.* The first two points are a result of repeatedly applying Lemma 4. By construction 64, for all  $i = 1, \dots$ , the  $i$ -th policy  $\pi_i$  is  $Q^{\pi_{i-1}}$ -greedy, implying

$$\pi_i >_V \pi_{i-1} \quad (65)$$

unless  $\pi_{i-1}$  is  $(Q^{\pi_{i-1}})$ -greedy, too, in which case

$$\pi_i =_V \pi_{i-1}. \quad (66)$$

This implies that the above iteration stops generating improved policies *if and only if* a greedy policy is generated. Let us denote the first time this happens by index  $N$ . Since  $n_{S,A} = |S| \cdot |A| < \infty$ , there are only a finite number of deterministic policies at our disposal. Since the sequence of policies

$$\pi_0 <_V \pi_1 < \dots < \pi_{N-1} < \pi_N \quad (67)$$

more generally implies

$$\pi_0 \neq \pi_1 \neq \dots \neq \pi_{N-1} \neq \pi_N, \quad (68)$$

we are essentially drawing without replacement from a finite pot of policies until  $\pi_N$  is reached and the sequence stagnates. Thus, the greedy policy break criteria has to be reached in a finite number of steps  $N \leq n_{S,A} < \infty$ .

What is left to show is that  $\pi_N$  is actually optimal. We will do this by a proof of contradiction.

If  $\pi_N$  is not optimal, there must be at least one other policy  $\pi_{imp} \in \Pi$  for which

$$\pi_{imp} >_V \pi_N \quad (69)$$

holds true. If  $\pi_{imp}$  itself is not greedy, we can apply the above iteration to  $\pi_0 = \pi_{imp}$  to obtain a second greedy policy,  $\pi_{imp, N_{imp}}$  satisfying

$$\pi_N <_V \pi_{imp, N_{imp}}. \quad (70)$$

By Corollary 6, this implies

$$\max_{a \in A} Q^{\pi_N}(s, a) < \max_{a \in A} Q^{\pi_{imp, N_{imp}}}(s, a) \quad (71)$$

for at least one  $s \in S$ , and, more generally,

$$Q^{\pi_N} \neq Q^{\pi_{imp, N_{imp}}}. \quad (72)$$

But both  $Q^{\pi_N}$  and  $Q^{\pi_{imp, N_{imp}}}$  are fixed points of  $B_{\mathbb{Q}}$ , which only has one fixed point by Proposition 4. This contradiction implies that  $\pi_N$  must have been optimal, and the proof is complete. □