

## Examples using Policy and Value Functions

(from *Sutton and Barto*)

Compiled by D. Gueorguiev 1/21/2024

### Gridworld

The cells of the grid correspond to the states of the environment. At each cell, four actions are possible: **north**, **south**, **east**, and **west** which deterministically cause the agent to move one cell in the respective direction on the grid. Actions that would take the agent off the grid leave its position unchanged, but also result in a reward of  $-1$ . Other actions result in a reward of 0, except those that move the agent out of the special states  $A$  and  $B$ . From state  $A$ , all four actions yield a reward of  $+10$  and take the agent to  $A'$ . From state  $B$ , all actions yield a reward of  $+5$  and take the agent to  $B'$ .

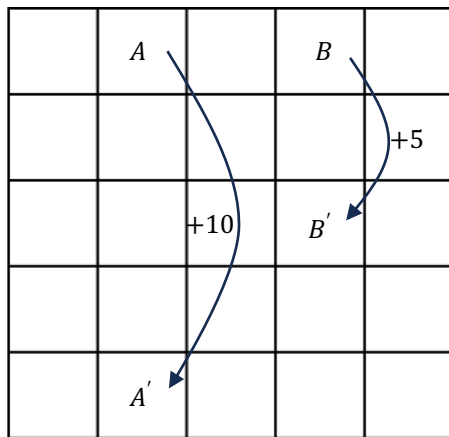


Figure 1: Gridworld Example

Suppose the agent selects all four actions with equal probability in all states. Figure 2 below shows the value function,  $v_\pi$ , for this policy, for the discounted reward case of  $\gamma = 0.9$ .

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Figure 2: State-value function for Gridworld

Now let us denote with  $s_i, i = 1..16$  and we set  $p = l/4, q = l \% 4$ .

Then  $s_i$  represents the state of the position  $(p, q), p, q \in [1,4]$  in the grid. Thus  $\mathcal{S} \equiv [1,16]$ .

We use the derived in the Appendix expression for (A.18)-(A.24):

$$p_{\pi}(s_j, r|s_i) \doteq \mathbb{E}[\pi|S_{t-1} = s_i] = \sum_a p(s_j, r|s_i, a) \cdot \pi(a|s_i) \quad (\text{A.18})$$

$$p_{\pi}(s_j|s_i) \doteq \sum_r p_{\pi}(s_j, r|s_i) \quad (\text{A.19})$$

$$r_{\pi}(s_j|s_i) \doteq \sum_r p_{\pi}(s_j, r|s_i) \cdot r \quad (\text{A.20})$$

For convenience we abbreviate:

$$p_{j,i} \doteq p_{\pi}(s_j|s_i), \quad r_i \doteq \sum_j r_{\pi}(s_j|s_i), \quad S \doteq |\mathcal{S}| \quad (\text{A.21})$$

Using (A.16)-(A.21) in (3) leads to :

$$(1 - \gamma p_{i,i})v_{\pi}(s_i) - \gamma \sum_{j \neq i} p_{j,i}v_{\pi}(s_j) = r_i \quad (\text{A.22})$$

(A.22) represents a linear system of equations with respect to the  $|\mathcal{S}|$  unknowns  $v_{\pi}(s_i), s_i \in \mathcal{S}$ .

Let us denote with  $\mathbf{v}_{\pi}$  the column vector of the  $S$  unknowns  $v_{\pi}(s_i), s_i \in \mathcal{S}$ .

We denote with  $\mathbf{A}$  the matrix formed by the elements and  $\mathbf{b}$  the vector as shown below:

$$\mathbf{A} = \begin{bmatrix} 1 - \gamma p_{11} & p_{21} & \dots & p_{S1} \\ p_{12} & 1 - \gamma p_{22} & \ddots & p_{S2} \\ & \vdots & \ddots & \vdots \\ p_{1S} & p_{2S} & \dots & 1 - \gamma p_{SS} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_S \end{bmatrix} \quad (\text{A.23})$$

Then (A.22) in matrix form:

$$\mathbf{A} \cdot \mathbf{v}_{\pi} = \mathbf{b} \quad (\text{A.24})$$

Let us assume random (equiprobable) policy  $\pi^{rand}(a|s_i)$  so that all  $a \in [north, south, east, west]$  are equally probable with probability  $1/4$ . The reward and the new state for every pair of action and state as a tuple  $(r, s')$  is given with the following table:

	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$	$s_8$	$s_9$	$s_{10}$	$s_{11}$	$s_{12}$	$s_{13}$	$s_{14}$	$s_{15}$	$s_{16}$
<i>north</i>	$(-1, s_1)$	$(-1, s_1)$	$(-1, s_1)$	$(-1, s_1)$												
<i>south</i>	$(0, s_5)$	$(0, s_6)$	$(0, s_7)$	$(0, s_8)$												
<i>east</i>	$(0, s_2)$															
<i>west</i>	$(-1, s_1)$															

Let us construct the function of the MDP dynamics  $p(s', r|s, a)$ .

## Appendix

Notation and Definitions from Sutton and Barto's RL book

**Distribution of the Dynamics of the MDP:** defined through the following 4 arguments function:

$$p(s', r|s, a) \doteq \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

which is the probability to get from state  $s$  to state  $s'$  with action  $a$  and with reward  $r$ .

**Distribution of the state-transition probabilities:** defined through the following 3 arguments function:

$$p(s'|s, a) \doteq \Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r|s, a)$$

**Markov Decision Process** (abbrev *MDP*): a 5-tuple  $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$  with

- $\mathcal{S}$  is a set of states (finite or infinite, discrete, or continuous)
- $\mathcal{A}$  is a set of actions (finite or infinite, discrete, or continuous)
- $p(s', r|s, a) \doteq \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$  is the function which describes the MDP dynamics i.e. probability to get from state  $s$  to state  $s'$  with action  $a$  and with reward  $r$ .
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  defines a *reward function*
- $\gamma \in [0, 1]$  is the discount factor which determines to what extent the focus is on the most recent rewards. with  $\gamma = 1$  there is no focus on the most recent rewards only.

Note: There is another equivalent definition of Markov process which uses the *state-transition probabilities distribution* represented by the three-argument function  $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ . With this definition the Markov Decision Process is defined as a 5-tuple  $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$ , where:

- $\mathcal{S}$  is a set of states (finite or infinite, discrete, or continuous)
- $\mathcal{A}$  is a set of actions (finite or infinite, discrete, or continuous)
- $p(s'|s, a) \doteq \Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\}$  is the distribution of the *state-transition probabilities* i.e. the probability to get from state  $s$  to state  $s'$  with action  $a$ .
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  defines a *reward function*
- $\gamma \in [0, 1]$  is the discount factor

Note 2: A more detailed definition of MDP involves specifying the initial state distribution  $d_0(s_0)$  and augments either of the MDP definitions as:

Markov Decision Process with specified *initial state* is a 6-tuple  $(\mathcal{S}, \mathcal{A}, T, r, d_0, \gamma)$ , where:

- $\mathcal{S}$  is a set of states (finite or infinite, discrete, or continuous)
- $\mathcal{A}$  is a set of actions (finite or infinite, discrete, or continuous)
- $p(s'|s, a) \doteq \Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\}$  is the probability to get from state  $s$  to state  $s'$  with action  $a$ .
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  defines a *reward function*
- $d_0(s_0)$  defines the initial state distribution
- $\gamma \in [0, 1]$  is the discount factor

**Learning Policy** (or just *Policy*): function  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  which represents mapping from states to probabilities of selecting each possible action.

If the agent is following policy  $\pi$  at time  $t$ , then  $\pi(a|s)$  is the probability that  $A_t = a$  if  $S_t = s$ . Note that  $\pi(a|s)$  is an ordinary function which defines a probability distribution over  $a \in \mathcal{A}(s)$  for each  $s \in \mathcal{S}$ .

We would like to modify the policy  $\pi$  with training or experience.

## State-Value and State-Action Functions

Let us assume that the current state is  $S_t$ , and actions are selected according to a stochastic policy  $\pi$ . Then we would like to derive an expression for the expectation of  $R_{t+1}$  in terms of  $\pi$  and  $p(s', r|s, a)$ .

Recall, the function  $p(s', r|s, a)$  defines the dynamics of the MDP and is given as:

$$p(s', r|s, a) \doteq \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\} \text{ for all } s', s \in \mathcal{S}, r \in \mathcal{R}, a \in \mathcal{A}(s) \quad (\text{A.1})$$

Then we can write:

$$\mathbb{E}_\pi[R_{t+1}|S_t = s] = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r] \quad (\text{A.2})$$

Here  $r$  denotes the reward of going from state  $s$  to state  $s'$  taking action  $a$  is given by MDP's  $R$  function:  $r = R(s, s', a)$ .

**State-Value Function for Policy  $\pi$**  (or simply *Value function*; aka *V function*): the value function of a state  $s$  under a policy  $\pi$ , denoted with  $v_\pi(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter. For MDPs, we can define  $v_\pi$  formally by

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t|S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s] \text{ for all } s \in \mathcal{S} \quad (\text{A.3})$$

where  $\mathbb{E}_\pi[\cdot]$  denotes the expected value of a random variable given that the agent follows policy  $\pi$ , and  $t$  is any time step. Note that the value of the terminal state, if any, is always zero.

**Action-Value Function for Policy  $\pi$**  (aka *Q function*):

We define the value of taking action  $a$  in state  $s$  under a policy  $\pi$ , denoted  $q_\pi(s, a)$ , as expected return starting from  $s$ , taking the action  $a$ , and thereafter following policy  $\pi$ :

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t|S_t = s, A_t = a] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a] \text{ for all } s \in \mathcal{S} \text{ and } a \in \mathcal{A}(s) \quad (\text{A.4})$$

Let us express  $v_\pi$  in terms of  $q_\pi$  and  $\pi$ . Given a state  $s$ , the state value function  $v_\pi(s)$ , given with (A.3), is equal to the expected cumulative return from that state given a distribution of actions  $\pi$ . The action value function  $q_\pi$  is the expectation of the return given state  $s$ , and taking action  $a$  as a starting point, and following policy  $\pi$  thereafter. Therefore, given a state  $s$ , the action-value function  $q_\pi$  is the weighted sum of the action-values over all relevant actions weighted by the policy weight:

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a) \quad (\text{A.5})$$

Given a state  $s$  and an action  $a$  let us express the action-value function  $q_\pi$  in terms of the state value function  $v_\pi$  and the function defining the MDP dynamics  $p(s', r|s, a)$ . Recall, given a state  $s$  and an action  $a$ , the action value function  $q_\pi$  is given by the mathematical expectation of the discounted future rewards i.e. return  $G_t$ . The return  $G_t$  is the discounted sequence of rewards after the time step  $t$  and it can be written as:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma G_{t+1} \quad (\text{A.6})$$

It is important to recognize that

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t|S_t = s, A_t = a] = \mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] + \gamma \mathbb{E}_\pi[G_{t+1}|S_t = s, A_t = a]. \quad (\text{A.7})$$

The first term on the right-hand side of (A.7) can be expressed as:

$$\mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] = \sum_{s'} \sum_r p(s', r|s, a) [r]. \quad (\text{A.8})$$

As before,  $r$  denotes the reward of going from state  $s$  to state  $s'$  taking action  $a$  is given by MDP's  $R$  function:  $r = R(s, s', a)$ .

The expectation in the second term on the right-hand side of (A.7) can be expressed as:

$$\mathbb{E}_\pi[G_{t+1}|S_t = s, A_t = a] = \sum_{s'} \sum_r p(s', r|s, a) v_\pi(s'). \quad (\text{A.9})$$

This is the expectation of the return starting at the next time step  $t + 1$  following the policy  $\pi$  given the current state  $s$  and the action  $a$ , chosen according to  $\pi$ .

Substituting (A.8) and (A.9) into (A.7) gives us:

$$q_{\pi}(s, a) = \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_{\pi}(s')]. \quad (\text{A.10})$$

Thus, the action-value function  $q_{\pi}$  given state  $s$  and action  $a$  following policy  $\pi$  is expressed as the sum of the next reward and discounted state-value weighted by probability distribution over the possible next states and next rewards from the given action  $a$  and state  $s$ .

## Bellman's Equations for State-Value and State-Action Functions

### Bellman's equation for state values $v$

The value functions satisfy recursive relationships, and this property of the former will prove quite useful. For any policy  $\pi$  and for any state  $s$ , the following consistency condition holds between the value of  $s$  and the value of its successor states. Starting with (A.6) applied to the definition of  $v_{\pi}(s)$ :

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (\text{A.11})$$

Using (5) the last equation becomes:

$$\begin{aligned} v_{\pi}(s) &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s']] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_{\pi}(s')] \quad \text{for all } s \in \mathcal{S} \end{aligned} \quad (\text{A.12})$$

where it is implicit that the actions,  $a$ , are taken from the set  $\mathcal{A}(s)$ , that the next states,  $s'$ , are taken from the set  $\mathcal{S}$ , and that the rewards,  $r$ , are taken from the set  $\mathcal{R}$ . Here the reward of going from state  $s$  to state  $s'$  taking action  $a$  is given by MDP's  $R$  function:  $r = R(s, s', a)$ . Note that the right-hand side of (A.12) is interpreted as an expected value obtained as a sum over the values of the triplet  $a, s'$ , and  $r$ . For each triplet  $(a, s', r)$  the quantity  $r + \gamma v_{\pi}(s')$  is weighed by its probability,  $\pi(a|s)p(s', r|s, a)$ .

Eq. (A.12) is known as the *Bellman equation* for  $v_{\pi}$ . It expresses a relationship between the value of a state and the values of its successor states.

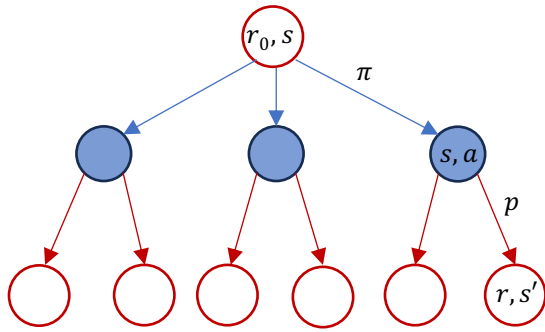


Figure A.1: Backup diagram for  $v_{\pi}$

This relationship is expressed by the *Backup diagram* shown on Figure A.1. Each open circle, which will be denoted as *reward-state node* so forth, colored in red represents a state and the reward, which is associated with this state. For instance, the root node shown on Figure 1 has associated reward  $r_0$  and state  $s$ . Each solid circle, colored in blue represents a state-action pair and will be denoted as *state-action node* so forth. The specific state on the rightmost state-action node is shown as  $(s, a)$ . Each directed blue edge connects state node with state-action

node and represents application of the policy  $\pi$  to the root reward-state node  $(r_0, s)$ . Each directed **red** edge emanating from a state-action node ends in a possible reward-state node corresponding to specific probable pair of reward  $r$  and new state  $s'$ . Thus, each directed **red** edge represents the application of the function  $p$  of the MDP dynamics. The Bellman equation (A.12) averages over all of the possibilities weighing each possibility represented by a path from the root of the Backup diagram on Figure A.1 to a leaf by its probability of occurring. It states that the value of the start state must equal the discounted value of the expected next state plus the reward expected along the way. The value function  $v_\pi$  is the unique solution to its Bellman equation. Various methods exist to compute exactly, approximate, or learn the value function.

### *Derivation of action value system of equations for discrete finite state*

Let us assume that we are dealing with *discrete finite state* – that is we have  $s_i \in \mathcal{S} \forall i \in [1, S]$ . Here we have denoted  $S = |\mathcal{S}|$ .

Then from (A.3) we have:

$$v_\pi(s_i) \doteq \mathbb{E}_\pi[G_t | S_t = s_i] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s_i] \text{ for all } s_i \in \mathcal{S} \quad (\text{A.13})$$

We also use (A.12) rewritten as:

$$v_\pi(s_i) = \sum_a \pi(a|s_i) \sum_{j \neq i, r} p(s_j, r|s_i, a) [r + \gamma v_\pi(s_j)] + \sum_a \pi(a|s_i) \sum_r p(s_i, r|s_i, a) [r + \gamma v_\pi(s_i)] \text{ for all } s_i, s_j \in \mathcal{S} \quad (\text{A.14})$$

hence

$$v_\pi(s_i) - \gamma \sum_a \pi(a|s_i) \sum_r p(s_i, r|s_i, a) \cdot v_\pi(s_i) - \gamma \sum_a \pi(a|s_i) \sum_{j \neq i, r} p(s_j, r|s_i, a) \cdot v_\pi(s_j) = \sum_a \pi(a|s_i) \sum_r p(s_i, r|s_i, a) \cdot r + \sum_a \pi(a|s_i) \sum_r p(s_i, r|s_i, a) \cdot r \quad (\text{A.15})$$

The left-hand side of (A2) can be rewritten as:

$$v_\pi(s_i) - \gamma \sum_r \sum_a p(s_i, r|s_i, a) \cdot \pi(a|s_i) v_\pi(s_i) - \gamma \sum_{j \neq i, r} \sum_a p(s_j, r|s_i, a) \cdot \pi(a|s_i) v_\pi(s_j) \quad (\text{A.16})$$

The right-hand side of (A.15) are rearranged as :

$$\sum_r \sum_a p(s_i, r|s_i, a) \cdot \pi(a|s_i) \cdot r + \sum_{j \neq i, r} \sum_a p(s_j, r|s_i, a) \cdot \pi(a|s_i) \cdot r \quad (\text{A.17})$$

we denote with  $p_\pi(s_j, r|s_i)$ ,  $p_\pi(s_j|s_i)$  and  $r_\pi(s_j|s_i)$  the following expressions:

$$p_\pi(s_j, r|s_i) \doteq \mathbb{E}[\pi | S_{t-1} = s_i] = \sum_a p(s_j, r|s_i, a) \cdot \pi(a|s_i) \quad (\text{A.18})$$

$$p_\pi(s_j|s_i) \doteq \sum_r p_\pi(s_j, r|s_i) \quad (\text{A.19})$$

$$r_\pi(s_j|s_i) \doteq \sum_r p_\pi(s_j, r|s_i) \cdot r \quad (\text{A.20})$$

For convenience we abbreviate:

$$p_{j,i} \doteq p_\pi(s_j|s_i), \quad r_i \doteq \sum_j r_\pi(s_j|s_i), \quad S \doteq |\mathcal{S}| \quad (\text{A.21})$$

Using (A.16)-(A.21) in (3) leads to :

$$(1 - \gamma p_{i,i}) v_\pi(s_i) - \gamma \sum_{j \neq i} p_{j,i} v_\pi(s_j) = r_i \quad (\text{A.22})$$

(A.22) represents a linear system of equations with respect to the  $|\mathcal{S}|$  unknowns  $v_\pi(s_i), s_i \in \mathcal{S}$ .

Let us denote with  $\mathbf{v}_\pi$  the column vector of the  $S$  unknowns  $v_\pi(s_i), s_i \in \mathcal{S}$ .

We denote with  $\mathbf{A}$  the matrix formed by the elements and  $\mathbf{b}$  the vector as shown below:

$$\mathbf{A} = \begin{bmatrix} 1 - \gamma p_{11} & p_{21} & \cdots & p_{s1} \\ p_{12} & 1 - \gamma p_{22} & \cdots & p_{s2} \\ \vdots & \vdots & \ddots & \vdots \\ p_{1s} & p_{2s} & \cdots & 1 - \gamma p_{ss} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_s \end{bmatrix} \quad (\text{A.23})$$

Then (A.22) in matrix form:

$$\mathbf{A} \cdot \mathbf{v}_\pi = \mathbf{b} \quad (\text{A.24})$$

### Bellman's equation for state-action values $q$

Let us derive a similar recursive relation with respect to the state-action value function. That is, we will find out what is the relation between the action value  $q_\pi(s, a)$  and that for the possible successors to the state-action pair  $(s, a)$ . The derivation follows from the Backup diagram shown on Figure A.2 below.

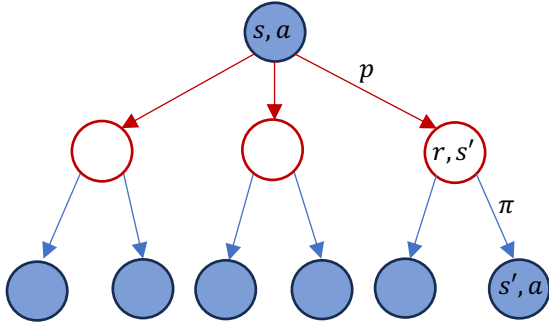


Figure A.2: Backup diagram for  $q_\pi$

From (A.7) we can write:

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \quad (\text{A.25})$$

The expectation of the reward on the right-hand side can be rewritten as:

$$\mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] = \sum_{s', r} p(s', r | s, a) [r + \gamma q_\pi(s', a)] \quad (\text{A.26})$$

Here using (A.7) again we denote with  $q_\pi(s', a')$  the expression for  $\mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a]$ .

Thus we get the Bellman's equation with respect  $q_\pi$ :

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma q_\pi(s', a)] \quad (\text{A.27})$$

### Expressing the current state values $v$ in terms of the next action values $q$

It is instructive to compare Eq (A.5) which we derived earlier with Eq (A.12) and Eq (A.27).  
Eq (A.5) deserves its own Backup diagram shown on Figure 3:

$$v_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a) \quad (\text{A.5})$$

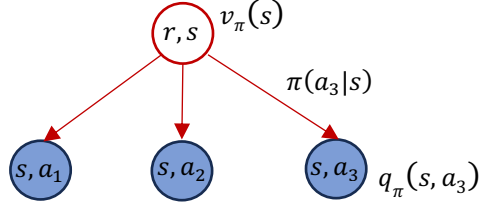


Figure A.3: Backup diagram for relation between  $v_{\pi}(s)$  and  $q_{\pi}(s, a)$

Eq (A.5) tells us how the value of a state depends on the values of the actions possible in that state and on how likely each action is to be taken under the current policy. The state value which corresponds to the state-value node at the root is obviously  $v_{\pi}(s)$  and the action values which corresponds to its children are  $q_{\pi}(s, a_i)$ ,  $i = 1..3$ . The probability with which each action  $a_i$  is taken is given by the policy i.e.  $\pi(a_i|s)$ ,  $i = 1..3$ .

### Expressing the current action values $q$ in terms of the next state values $v$

The value of an action,  $q_{\pi}(s, a)$ , depends on the expected next reward and the expected sum of the remaining rewards. Expressed as a Backup diagram we arrive at Figure A.4 shown below.

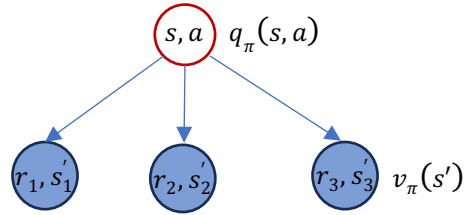


Figure A.4: Backup diagram expressing the dependence of the current action value on the expected next reward-state values.

Formally expressed this relation becomes:

From Eq (A.7) we have

$$q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi}[R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_t = s, A_t = a]$$

Clearly,  $\mathbb{E}_{\pi}[R_{t+1} | S_t = s, A_t = a] = \sum_{s', r} p(s', r | s, a) \cdot r$  where  $r = R(s, s', a)$

Eq. (A.9) states that the expectation in the second term of (A.7) can be written as:

$$\mathbb{E}_{\pi}[G_{t+1} | S_t = s, A_t = a] = \sum_{s', r} p(s', r | s, a) \cdot v_{\pi}(s')$$

Combining the last two results we obtain:

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(s') | S_t = s, A_t = a] = \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \quad (\text{A.28})$$



## Bibliography

[Reinforcement Learning, Richard S. Sutton, Andrew G. Barto, second edition, 2020](#)