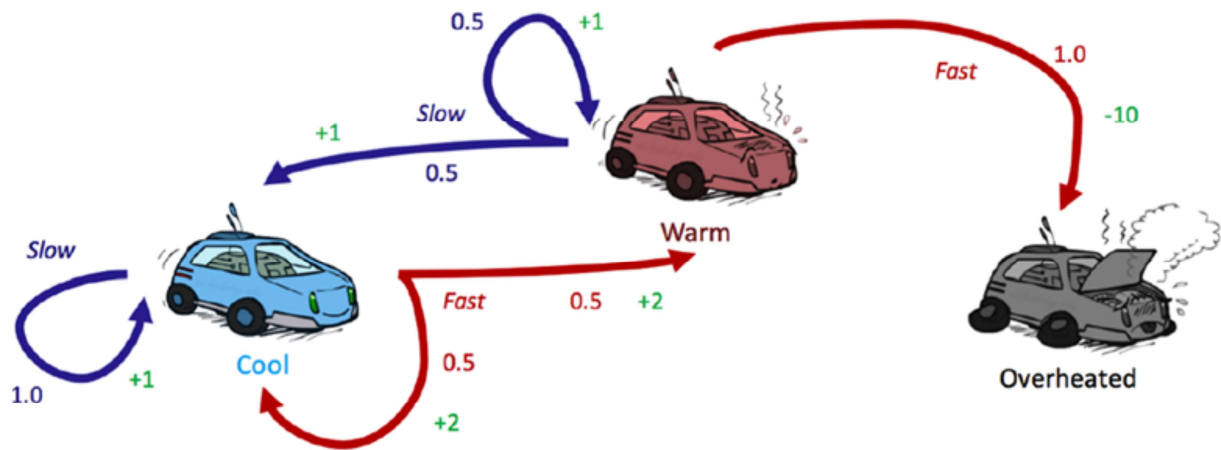




# Complex Decisions

Chapter 17



$$U(s) = \max_{a \in A(s)} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma U(s')] ]$$

States={cool, warm, overheated}

Actions={fast, slow}

$\gamma=0.5$

Initially:  $U_0(\text{cool})=0$ ,  $U_0(\text{warm})=0$ ,  $U_0(\text{overheated})=0$

Run two iterations of the value iteration algorithm.

$$\begin{aligned}
 U_{t+1}(\text{cool}) &= \max( \\
 &\quad P(\text{cool} | \text{cool}, \text{slow}) * (R(\text{cool}, \text{slow}, \text{cool}) + \gamma U_t(\text{cool})), \\
 &\quad P(\text{cool} | \text{cool}, \text{fast}) * (R(\text{cool}, \text{fast}, \text{cool}) + \gamma U_t(\text{cool})) + P(\text{warm} | \text{cool}, \text{fast}) * \\
 &\quad (R(\text{cool}, \text{fast}, \text{warm}) + \gamma U_t(\text{warm})) \\
 &= \max(1 * (1 + 0.5U_t(\text{cool})), 0.5 * (2 + 0.5U_t(\text{cool})) + 0.5 * (2 + 0.5U_t(\text{warm}))) = \max(1 + \\
 &\quad 0.5U_t(\text{cool}), 2 + 0.25U_t(\text{cool}) + 0.25U_t(\text{warm})) \\
 U_{t+1}(\text{warm}) &= \max( \\
 &\quad P(\text{cool} | \text{warm}, \text{slow}) * (R(\text{warm}, \text{slow}, \text{cool}) + \gamma U_t(\text{cool})) + P(\text{warm} | \text{warm}, \text{slow}) \\
 &\quad * (R(\text{warm}, \text{slow}, \text{warm}) + \gamma U_t(\text{warm})), \\
 &\quad P(\text{overheated} | \text{warm}, \text{fast}) * (R(\text{warm}, \text{fast}, \text{overheated}) + \gamma U_t(\text{overheated})) \\
 &= \max(0.5 * (1 + 0.5U_t(\text{cool})) + 0.5 * (1 + 0.5U_t(\text{warm})), 1 * (-10 + 0.5U_t(\text{overheated}))) = \\
 &\quad \max(1 + 0.25U_t(\text{cool}) + 0.25U_t(\text{warm}), -10 + 0.5U_t(\text{overheated})) \\
 U_{t+1}(\text{overheated}) &= 0
 \end{aligned}$$

Iteration 1:

$$U_1(\text{cool}) = \max(1 + 0.5U_0(\text{cool}), 2 + 0.25U_0(\text{cool}) + 0.25U_0(\text{warm})) = \max(1 + 0.5*0, 2 + 0.25*0 + 0.25*0) = \max(1, 2) = 2$$

$$U_1(\text{warm}) = \max(1 + 0.25U_0(\text{cool}) + 0.25U_0(\text{warm}), -10 + 0.5U_0(\text{overheated})) = \max(1 + 0.25*0 + 0.25*0, -10 + 0.5*0) = \max(1, -10) = 1$$

$$U_1(\text{overheated}) = 0$$

Iteration 2:

$$U_2(\text{cool}) = \max(1 + 0.5U_1(\text{cool}), 2 + 0.25U_1(\text{cool}) + 0.25U_1(\text{warm})) = \max(1 + 0.5*2, 2 + 0.25*2 + 0.25*1) = \max(2, 2.75) = 2.75$$

$$U_2(\text{warm}) = \max(1 + 0.25U_1(\text{cool}) + 0.25U_1(\text{warm}), -10 + 0.5U_1(\text{overheated})) = \max(1 + 0.25*2 + 0.25*1, -10 + 0.5*0) = \max(1.75, -10) = 1.75$$

$$U_2(\text{overheated}) = 0$$

Policy extraction:

$$\begin{aligned}\pi(\text{cool}) &= \operatorname{argmax}(\text{slow: } 1 + 0.5U_1(\text{cool}), \text{ fast: } 2 + 0.25U_1(\text{cool}) + 0.25U_1(\text{warm})) = \\ &\operatorname{argmax}(\text{slow: } 1 + 0.5 \cdot 2.75, \text{ fast: } 2 + 0.25 \cdot 2.75 + 0.25 \cdot 1.75) = \operatorname{argmax}(\text{slow: } 2.375, \text{ fast: } \\ &3.125) = \text{fast}\end{aligned}$$

$$\begin{aligned}\pi(\text{warm}) &= \operatorname{argmax}(\text{slow: } 1 + 0.25U_1(\text{cool}) + 0.25U_1(\text{warm}), \text{ fast: } -10 + 0.5U_1(\text{overheated})) \\ &= \operatorname{argmax}(\text{slow: } 1 + 0.25 \cdot 2.75 + 0.25 \cdot 1.75, \text{ fast: } -10 + 0.5 \cdot 0) = \operatorname{argmax}(\text{slow: } 1.375, \text{ fast: } - \\ &-10) = \text{slow}\end{aligned}$$

**17.4** Sometimes MDPs are formulated with a reward function  $R(s, a)$  that depends on the action taken or with a reward function  $R(s, a, s')$  that also depends on the outcome state.

- a. Write the Bellman equations for these formulations.
- b. Show how an MDP with reward function  $R(s, a, s')$  can be transformed into a different MDP with reward function  $R(s, a)$ , such that optimal policies in the new MDP correspond exactly to optimal policies in the original MDP.
- c. Now do the same to convert MDPs with  $R(s, a)$  into MDPs with  $R(s)$ .

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$$U(s) = \max_{a \in A(s)} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) U(s')]$$

$$U(s) = \max_{a \in A(s)} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma U(s')]$$

a. The key here is to get the max and summation in the right place.

For  $R(s)$ , we use the first formula above

For  $R(s, a)$ , we use the second formula above

For  $R(s, a, s')$ , we use the third formula above.

b. There are a variety of solutions here. One is to create a “pre-state”  $\text{pre}(s, a, s')$  for every  $s, a, s'$ , such that executing  $a$  in  $s$  leads not to  $s'$  but to  $\text{pre}(s, a, s')$ . In this state is encoded the fact that the agent came from  $s$  and did  $a$  to get here. From the pre-state, there is just one action  $b$  that always leads to  $s'$ . Let the new MDP have transition  $P'$ , reward  $R'$ , and discount  $\gamma'$ . Then

$$P'(\text{pre}(s, a, s') | s, a) = P(s' | s, a)$$

$$P'(s' | \text{pre}(s, a, s'), b) = 1$$

$$R'(s, a) = 0$$

$$R'(\text{pre}(s, a, s'), b) = R(s, a, s') / \sqrt{\gamma}$$

$$\gamma' = \sqrt{\gamma}$$

c. In keeping with the idea of part (b), we can create states  $\text{post}(s, a)$  for every  $s, a$ , such that

$$P'(\text{post}(s, a, s') | s, a) = 1$$

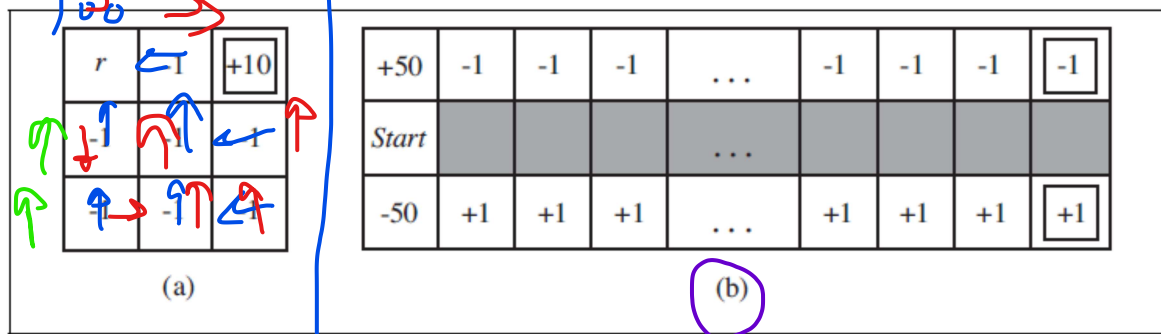
$$P'(s' | \text{post}(s, a, s'), b) = P(s' | s, a)$$

$$R'(s) = 0$$

$$R'(\text{post}(s, a, s')) = R(s, a) / \sqrt{\gamma}$$

$$\gamma' = \sqrt{\gamma}$$

+ } same as 100  
0



**17.8** Consider the  $3 \times 3$  world shown in Figure 17.14(a). The transition model is the same as in the  $4 \times 3$  Figure 17.1: 80% of the time the agent goes in the direction it selects; the rest of the time it moves at right angles to the intended direction.

Implement value iteration for this world for each value of  $r$  below. Use discounted rewards with a discount factor of 0.99. Show the policy obtained in each case. Explain intuitively why the value of  $r$  leads to each policy.

- ✓ a.  $r = 100$
- (b)  $r = -3$
- c.  $r = 0$
- d.  $r = +3$

Instead of actually implementing value iteration, try to predict the policy by intuition

- a.
- LL\_
- ULD
- ULL

Instead of seeking the goal, the agent would prefer to go the  $r=100$  cell. Since the discount factor is high, the agent wouldn't mind paying some losses to reach a larger reward later on. So the agent will avoid the terminal cell (+10) and will do whatever it can to reach the +100 cell. At the cell below the +10, it would take the action "Down" to avoid entering the terminal cell by accident. At the cells in the middle column, it will always pick "Left" to make sure it never goes right by accident.

At some cells, there are multiple optimal actions. For example, at the "+100" cell, it can pick "Up" or "Left" and the utility shouldn't differ.

- b.
- RR\_
- RRU
- RRU

The agent will try to avoid the "-3" cell and try to reach the terminal state "+10" as soon as possible.

- c.
- RR\_
- UUU
- UUU

The agent will try to reach the terminal state "+10" as soon as possible and it will prefer the path going through the "0"-reward cell (if there is no other path that is shorter).

d.

LL\_

ULD

ULL

This is similar to the first one. It may seem surprising that the agent will prefer the “+3” cell over the terminal “+10” cell but remember that the “+3” cell is not a terminal state. Thus the agent can keep getting rewards from it as long as it stays on it. So the agent can gain much more rewards by staying on the “+3” cell as long as possible than what it would get from reach the “+10” state which would return +10 once then end the episode. This is only true because the discount factor is high. If it was low, the agent would prefer getting the “+10” early than waiting to get multiple “+3”s in the future.

<i>r</i>	-1	<span style="border: 1px solid black;">+10</span>
-1	-1	-1
-1	-1	-1

(a)

+50	-1	-1	-1	...	-1	-1	-1	<span style="border: 1px solid black;">-1</span>
<i>Start</i>				...				
-50	+1	+1	+1	...	+1	+1	+1	<span style="border: 1px solid black;">+1</span>

(b)

**17.9** Consider the  $101 \times 3$  world shown in Figure 17.14(b). In the start state the agent has a choice of two deterministic actions, *Up* or *Down*, but in the other states the agent has one deterministic action, *Right*. Assuming a discounted reward function, for what values of the discount  $\gamma$  should the agent choose *Up* and for which *Down*? Compute the utility of each action as a function of  $\gamma$ . (Note that this simple example actually reflects many real-world situations in which one must weigh the value of an immediate action versus the potential continual long-term consequences, such as choosing to dump pollutants into a lake.)

$$U((1,1)) = -50 + \gamma^1 + \gamma^2 + \gamma^3 + \dots + \gamma^{100} = -50 + \gamma(1 - \gamma^{100})/(1 - \gamma)$$

$$U((1,3)) = 50 - \gamma^1 - \gamma^2 - \gamma^3 - \dots - \gamma^{100} = 50 - \gamma(1 - \gamma^{100})/(1 - \gamma)$$

To choose up, it must be  $U((1,1)) < U((1,3))$

$$-50 + \gamma(1 - \gamma^{100})/(1 - \gamma) < 50 - \gamma(1 - \gamma^{100})/(1 - \gamma)$$

$$2\gamma(1 - \gamma^{100})/(1 - \gamma) < 2(50)$$

$$\gamma(1 - \gamma^{100}) < 50(1 - \gamma) \quad \text{since } 1 - \gamma > 0$$

This can be solved using numerical methods.

$$\gamma < 0.9843976692$$

So if  $\gamma < 0.9843976692$ , we choose up

Otherwise, we choose down



**17.10** Consider an undiscounted MDP having three states, (1, 2, 3), with rewards  $-1$ ,  $-2$ ,  $0$ , respectively. State 3 is a terminal state. In states 1 and 2 there are two possible actions:  $a$  and  $b$ . The transition model is as follows:

- In state 1, action  $a$  moves the agent to state 2 with probability 0.8 and makes the agent stay put with probability 0.2.
- In state 2, action  $a$  moves the agent to state 1 with probability 0.8 and makes the agent stay put with probability 0.2.
- In either state 1 or state 2, action  $b$  moves the agent to state 3 with probability 0.1 and makes the agent stay put with probability 0.9.

Answer the following questions:

- What can be determined *qualitatively* about the optimal policy in states 1 and 2?
- Apply policy iteration, showing each step in full, to determine the optimal policy and the values of states 1 and 2. Assume that the initial policy has action  $b$  in both states.
- What happens to policy iteration if the initial policy has action  $a$  in both states? Does discounting help? Does the optimal policy depend on the discount factor?

**Extra: Apply value iteration for 3 iterations then extract the policy.**

a. The policy should be (b, a, \_)

Because, the agent will try reach the terminal state (state 3) as soon as possible to stop getting penalties from being in state 1 or 2.

The only action to reach state 3 is “b” which only has a probability of 0.1 to actually move the agent to state 3. So the agent should prefer to go to state 1 first since it has a lower penalty then try to go to state 3 from there. Therefore, the action to pick in state 2 is “a” while the action to pick in state 1 is “b”.

b. Since  $\pi(1) = b$  and  $\pi(2) = b$ , the values for this policy can be computed from using the following equation:

$$U(1) = -1 + 0.1 \cdot U(3) + 0.9 \cdot U(1)$$

$$U(2) = -2 + 0.1 \cdot U(3) + 0.9 \cdot U(2)$$

$$U(3) = 0$$

By solving this equation, we get  $U(1) = -10$ ,  $U(2) = -20$ ,  $U(3) = 0$

Now, we compute the new policy

$$\pi(1) = \operatorname{argmax}(a: 0.8 \cdot U(2) + 0.2 \cdot U(1), b: 0.1 \cdot U(3) + 0.9 \cdot U(1)) = \operatorname{argmax}(a: 0.8 \cdot -20 + 0.2 \cdot -10, b: 0.1 \cdot 0 + 0.9 \cdot -10) = \operatorname{argmax}(a: -18, b: -9) = b$$

$$\pi(2) = \operatorname{argmax}(a: 0.8 \cdot U(1) + 0.2 \cdot U(2), b: 0.1 \cdot U(3) + 0.9 \cdot U(2)) = \operatorname{argmax}(a: 0.8 \cdot -10 + 0.2 \cdot -20, b: 0.1 \cdot 0 + 0.9 \cdot -20) = \operatorname{argmax}(a: -12, b: -18) = a$$

Then we get the new equations:

$$U(1) = -1 + 0.1 \cdot U(3) + 0.9 \cdot U(1)$$

$$U(2) = -2 + 0.8 \cdot U(1) + 0.2 \cdot U(2)$$

$$U(3) = 0$$

So  $U(1) = -10$ ,  $U(2) = -12.5$ ,  $U(3) = 0$

Now, we compute the new policy

$$\pi(1) = \operatorname{argmax}(a: 0.8 \cdot U(2) + 0.2 \cdot U(1), b: 0.1 \cdot U(3) + 0.9 \cdot U(1)) = \operatorname{argmax}(a: 0.8 \cdot -12.5 + 0.2 \cdot -10, b: 0.1 \cdot 0 + 0.9 \cdot -10) = \operatorname{argmax}(a: -12, b: -9) = b$$

10, b:  $0.1 \cdot 0 + 0.9 \cdot -10$ ) =  $\text{argmax}(a: -12, b: -9) = b$

$\pi(2) = \text{argmax}(a: 0.8 \cdot U(1) + 0.2 \cdot U(2), b: 0.1 \cdot U(3) + 0.9 \cdot U(2)) = \text{argmax}(a: 0.8 \cdot -10 + 0.2 \cdot -12.5, b: 0.1 \cdot 0 + 0.9 \cdot -12.5) = \text{argmax}(a: -10.5, b: -11.25) = a$

Since the policy did not change from the last iteration, we stop now.

The optimal policy is (b, a, \_).

c. Since  $\pi(1) = a$  and  $\pi(2) = a$ , the values for this policy can be computed from using the following equation:

$$U(1) = -1 + 0.8 \cdot U(2) + 0.2 \cdot U(1)$$

$$U(2) = -2 + 0.8 \cdot U(1) + 0.2 \cdot U(2)$$

$$U(3) = 0$$

These equations has no solution.

If we add a discount factor  $\gamma$ , the equations will become

$$U(1) = -1 + \gamma(0.8 \cdot U(2) + 0.2 \cdot U(1))$$

$$U(2) = -2 + \gamma(0.8 \cdot U(1) + 0.2 \cdot U(2))$$

$$U(3) = 0$$

Which are solvable if  $\gamma < 1$

And yes, the optimal policy depends on the discount factor.

Extra:

Assume that initially,  $U(1)=U(2)=U(3)=0$

Iteration 1:

$$U(1) = -1 + \max(0.8 \cdot U(2) + 0.2 \cdot U(1), 0.1 \cdot U(3) + 0.9 \cdot U(1)) = -1 + \max(0.8 \cdot 0 + 0.2 \cdot 0, 0.1 \cdot 0 + 0.9 \cdot 0) = -1 + \max(0, 0) = -1 + 0 = -1$$

$$U(2) = -2 + \max(0.8 \cdot U(1) + 0.2 \cdot U(2), 0.1 \cdot U(3) + 0.9 \cdot U(2)) = -2 + \max(0.8 \cdot 0 + 0.2 \cdot 0, 0.1 \cdot 0 + 0.9 \cdot 0) = -2 + \max(0, 0) = -2 + 0 = -2$$

$$U(3) = 0$$

Iteration 2:

$$U(1) = -1 + \max(0.8 \cdot U(2) + 0.2 \cdot U(1), 0.1 \cdot U(3) + 0.9 \cdot U(1)) = -1 + \max(0.8 \cdot -2 + 0.2 \cdot -1, 0.1 \cdot 0 + 0.9 \cdot -1) = -1 + \max(-1.8, -0.9) = -1 + -0.9 = -1.9$$

$$U(2) = -2 + \max(0.8 \cdot U(1) + 0.2 \cdot U(2), 0.1 \cdot U(3) + 0.9 \cdot U(2)) = -2 + \max(0.8 \cdot -1.9 + 0.2 \cdot -2, 0.1 \cdot 0 + 0.9 \cdot -2) = -2 + \max(-1.2, -1.8) = -2 + -1.2 = -3.2$$

$$U(3) = 0$$

Iteration 3:

$$U(1) = -1 + \max(0.8 \cdot U(2) + 0.2 \cdot U(1), 0.1 \cdot U(3) + 0.9 \cdot U(1)) = -1 + \max(0.8 \cdot -3.2 + 0.2 \cdot -1.9, 0.1 \cdot 0 + 0.9 \cdot -1.9) = -1 + \max(-2.94, -1.71) = -1 + -1.71 = -2.71$$

$$U(2) = -2 + \max(0.8 \cdot U(1) + 0.2 \cdot U(2), 0.1 \cdot U(3) + 0.9 \cdot U(2)) = -2 + \max(0.8 \cdot -2.71 + 0.2 \cdot -3.2, 0.1 \cdot 0 + 0.9 \cdot -3.2) = -2 + \max(-2.16, -2.38) = -2 + -2.16 = -4.16$$

$$U(3) = 0$$

Now we extract the policy:

$$\pi(1) = \text{argmax}(a: 0.8 \cdot U(2) + 0.2 \cdot U(1), b: 0.1 \cdot U(3) + 0.9 \cdot U(1)) = \text{argmax}(a: 0.8 \cdot -4.16 + 0.2 \cdot -2.71, b: 0.1 \cdot 0 + 0.9 \cdot -2.71) = \text{argmax}(a: -3.37, b: -2.439) = b$$

$$\pi(2) = \text{argmax}(a: 0.8 \cdot U(1) + 0.2 \cdot U(2), b: 0.1 \cdot U(3) + 0.9 \cdot U(2)) = \text{argmax}(a: 0.8 \cdot -2.71 + 0.2 \cdot -4.16, b: 0.1 \cdot 0 + 0.9 \cdot -4.16) = \text{argmax}(a: -3, b: -3.744) = a$$

Handwritten notes in purple ink:

- 6/1/24
- Done