

## Reward Shaping

Lecturer: Pieter Abbeel

Scribe: Pål From

# 1 Algorithm Review

## 1.1 Q-learning

The Q-learning algorithm can be summarized as

$$\begin{aligned}
 &\text{initialize } s_0 \\
 &\text{for } t = 1, 2, 3 \dots \\
 &\quad \text{choose an action } a_t, \text{ let's say } \epsilon\text{-greedy w.r.t. } Q \\
 &\quad \text{execute } a_t \\
 &\quad Q(s_t, a_t) = (1 - \alpha_t)Q(s_t, a_t) + \alpha_t [R(s_t, a_t, s_{t+1}) + \gamma \max_a Q(s_{t+1}, a)]
 \end{aligned} \tag{1}$$

We can for example choose  $\alpha_t = \frac{1}{t+1}$ . Although we have not shown it, we know this algorithm converges as long as we visit all the state with probability  $P(\cdot) > 0$ .

## 1.2 Value iteration

Value iteration is summarized as

$$\begin{aligned}
 &\text{iterate} \\
 &\quad \forall s \quad \bar{V}(s) = \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma V(s')] \\
 &\quad \forall s \quad V = \bar{V}
 \end{aligned} \tag{2}$$

which we normally write as  $V = TV$ .

# 2 Example

We want to find the optimal policy to move from the start to the goal state.

	S				G
Reward for entering state	0	0	0	0	1
State	1	2	3	4	5

(3)

$$\begin{aligned}
 S &- \text{ Start} \\
 G &- \text{ Goal}
 \end{aligned} \tag{4}$$

We have two actions

$$\begin{aligned}
 \hat{L} &- \text{ move to the state on the left} \\
 \hat{R} &- \text{ move to the state on the right}
 \end{aligned} \tag{5}$$

## 2.1 Q-learning

We apply Q-learning to the problem: We start by initializing  $s_0 = 1$  and  $a = \hat{R}$

$$Q(1, \hat{R}) = (1 - \alpha) \underbrace{Q(1, \hat{R})}_0 + \alpha \left[ \underbrace{R(1, \hat{R}, 2)}_0 + \gamma \max_a \underbrace{Q(2, a)}_0 \right] = 0 \quad (6)$$

we now have  $s_1 = 2$  and choose  $a = \hat{L}$

$$Q(2, \hat{L}) = (1 - \alpha) \underbrace{Q(2, \hat{L})}_0 + \alpha \left[ \underbrace{R(2, \hat{L}, 1)}_0 + \gamma \max_a \underbrace{Q(1, a)}_0 \right] = 0 \quad (7)$$

We see that until we coincidentally happen to hit state 5  $Q$  will always be zero. Very little information is added to the Q-learning algorithm before this happens as there is nothing guiding us towards the right state.

We now propose the alternative reward function

	S			G	
Reward for entering state	0	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{3}{4}$	1
State	1	2	3	4	5

(8)

We now get

$$Q(1, \hat{R}) = \alpha \cdot \gamma \cdot \frac{1}{4} \quad (9)$$

$$Q(2, \hat{L}) = 0 \quad (10)$$

Then at time 3 we are in state 1 and have

$$Q(1, \hat{R}) = \alpha \cdot \gamma \cdot \frac{1}{4}, \quad Q(1, \hat{L}) = 0 \quad (11)$$

We see that we can shape the reward to guide us in the right direction.

## 2.2 Value Iteration

State	0	1	2	...
Value Function	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \gamma \cdot 1 \\ \gamma \cdot 1 + 1 \\ \gamma \cdot 1 + 1 \end{bmatrix}$	...

(12)

We do this propagating from back to front. After four iterations we have found the optimal policy, although we have not yet converged to the value function.

## 3 Reward shaping without changing the optimal policy

So far the approach has been very heuristic. We will now look into how we can shape the reward function without changing the relative optimality of policies.

We start by looking at a bad example: let's say we want an agent to reach a goal state for which it has to climb over three mountains to get there. The original reward function has a zero reward everywhere, and a positive reward at the goal state (which is beyond the three mountains). We could imagine changing the

reward by adding a positive reward for making progress towards the goal by adding a positive reward when the agent reaches the top of each mountain. Indeed, this would favor some policies that move the agent towards the goal. However, it would also favor the following policy: climb to the top of the first mountain, take a few steps back, go back to the top, and keep repeating ... In fact, depending on the discounting and the exact reward function, the latter could even be the optimal policy.

Intuitively, the reward shaping in the example fails because the agent gets rewarded every time they reach the top of the mountain, independent of whether they already reached the top before. A natural solution is to, indeed reward the agent for reaching the top, but also *penalize* the agent for moving away from the goal/top such that cyclic behaviour results in a zero reward. This suggests shaping the reward function by using a potential function  $\phi$ :

$$\bar{R}(s, a, s') = R(s, a, s') + F(s, a, s') \quad (13)$$

where

$$F(s, a, s') = \phi(s') - \phi(s). \quad (14)$$

For now, we assume:

1.  $\exists$  “absorbing” state  $s_F$
2. No discounting

The first assumption means that for all policies, we will end up in state  $s_t = s_F$  with probability  $P = 1$ , i.e.

$$\forall \pi \quad \text{eventually } s_t = s_F. \quad (15)$$

We can remove the restriction on the discounting at a later stage.

Now consider the rewards accumulated during one episode:

$$\sum_{t=0}^{\tau-1} [R(s_t, a_t, s_{t+1}) + \phi(s_{t+1}) - \phi(s_t)] = \left[ \sum_{t=0}^{\tau-1} R(s_t, a_t, s_{t+1}) \right] + \phi(s_F) - \phi(s_0) \quad (16)$$

Let  $\bar{M}$  denote the MDP identical to the original MDP  $M$ , except for the reward function: in  $\bar{M}$  we have the shaped reward function. Then the above can be written as:

$$\forall \pi : V_{\bar{M}}^{\pi}(s) = V_M^{\pi}(s) + \phi(s_{\tau}) - \phi(s_0) \quad (17)$$

This implies that for all policies  $\pi_1, \pi_2$  we have that the ordering of their value is preserved between  $M$  and  $\bar{M}$ :

$$V_M^{\pi_1}(s) > V_M^{\pi_2}(s) \implies V_{\bar{M}}^{\pi_1}(s) > V_{\bar{M}}^{\pi_2}(s) \quad (18)$$

Hence reward shaping based upon differencing a potential function has the desired property of keeping the optimality ordering of policies invariant.

### 3.1 Infinite horizon

We re-write (16) with the discount factor

$$\sum_{t=0}^{\tau-1} \gamma^t [R(s_t, a_t, s_{t+1}) + \phi(s_{t+1}) - \phi(s_t)] \quad (19)$$

We write out

$$\sum_{t=0}^{\tau-1} \gamma^t [\phi(s_{t+1}) - \phi(s_t)] = (\phi(s_1) - \phi(s_0)) + (\gamma\phi(s_2) - \gamma\phi(s_1)) + (\gamma^2\phi(s_3) - \gamma^2\phi(s_2)) \dots \quad (20)$$

We see that if we multiply the first part with  $\gamma$ , we get

$$\begin{aligned} \sum_{t=0}^{\tau-1} \gamma^t [\gamma\phi(s_{t+1}) - \phi(s_t)] &= \gamma\phi(s_1) - \phi(s_0) + \gamma^2\phi(s_2) - \gamma\phi(s_1) + \gamma^3\phi(s_3) - \gamma^2\phi(s_2) \dots \\ &= -\phi(s_0) + \underbrace{(\gamma\phi(s_1) - \gamma\phi(s_1))}_0 + \underbrace{(\gamma^2\phi(s_2) - \gamma^2\phi(s_2))}_0 + \gamma^{\tau-1}\phi(s_\tau) \dots \end{aligned} \quad (21)$$

So we choose

$$F(s, a, s') = \gamma\phi(s') - \phi(s). \quad (22)$$

**Proposition 1** *Reward shaping with the function  $F(s, a, s') = \gamma\phi(s') - \phi(s)$  leaves the optimality ordering of policies invariant.*

**Proof** Let  $M$  be the original MDP, and let  $\bar{M}$  be identical except for the reward function being shaped by  $F(s, a, s') = \gamma\phi(s') - \phi(s)$ .

$$Q_M^*(s, a) = \sum_{s'} P(s' | s, a) \left[ R(s, a, s') + \gamma \max_a Q_M^*(s, a) \right] \quad (23)$$

we add the shaping (by adding and subtracting the same terms)

$$\underbrace{Q_M^*(s, a) - \phi(s)}_{Q_M^*(s, a)} = \sum_{s'} P(s' | s, a) \left[ \underbrace{R(s, a, s') + \gamma\phi(s') - \phi(s)}_{\text{Shaped Reward Function } R_{\bar{M}}(s, a, s')} + \gamma \max_a \underbrace{(Q_M^*(s, a) - \phi(s'))}_{Q_{\bar{M}}^*(s', a)} \right] \quad (24)$$

For  $\bar{M}$  we get

$$Q_{\bar{M}}^*(s, a) = \sum_{s'} P(s' | s, a) \left[ R_{\bar{M}}(s, a, s') + \gamma \max_a Q_{\bar{M}}^*(s', a) \right] \quad (25)$$

We note that  $Q_{\bar{M}}^*(s, a) = Q_M^*(s, a) - \phi(s)$  satisfies the Bellman equation for the reward  $R_{\bar{M}}(s, a, s') = R(s, a, s') + \gamma\phi(s') - \phi(s)$ . Thus we have

$$\begin{aligned} Q_{\bar{M}}^*(s, a) &= Q_M^*(s, a) - \phi(s) \\ \arg \max_a Q_{\bar{M}}^*(s, a) &= \arg \max_a Q_M^*(s, a) \end{aligned} \quad (26)$$

as we wanted and the optimal policy is preserved.

Now consider a special MDP, in which only 1 action is available in each state, namely the action prescribed by a policy  $\pi$ . Then the same reasoning as above goes through, and we obtain:

$$Q_{\bar{M}}^\pi(s, a) = Q_M^\pi(s, a) - \phi(s). \quad (27)$$

This holds true for any policy  $\pi$ , hence for any pair of policies  $\pi_1, \pi_2$  we have:

$$Q_{\bar{M}}^{\pi_1}(s, \pi_1(s)) \geq Q_{\bar{M}}^{\pi_2}(s, \pi_2(s)) \quad (28)$$

implies:

$$Q_M^{\pi_1}(s, \pi_1(s)) - \phi(s) \geq Q_M^{\pi_2}(s, \pi_2(s)) - \phi(s) \quad (29)$$

hence:

$$Q_M^{\pi_1}(s, \pi_1(s)) \geq Q_M^{\pi_2}(s, \pi_2(s)). \quad (30)$$

.667em

## 4 What is ideal shaping?

Consider value iteration, where we iteratively compute the value function as follows:

for  $k = 0, 1, \dots$

$$\forall s \quad V_{k+1}(s) = \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma V_k(s')]$$

With reward shaping we have:

$$V_{k+1}(s) = \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma \phi(s') - \phi(s) + \gamma V_k(s')]$$

Assume we choose  $V_0 = 0$ , and  $\phi = V^*$ , then we obtain:

$$V_1(s) = \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma V^*(s') - V^*(s) + 0] = V^*(s) - V^*(s) = 0$$

And similarly, for all  $k > 0$ , we have  $V_k(s) = 0$  for all states  $s$ . Note we also have:

$$\arg \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma V^*(s') - V^*(s) + V_k(s')] = \arg \max_a \sum_{s'} P(s'|a, s) [R(s, a, s') + \gamma V^*(s')] = \pi^*(s)$$

Hence, from the first iteration onwards, we have converged to the optimal value function and optimal policy.

This also suggests there is a close connection between initializing the value function and potential shaping.