

Blatt 4

Q-learning

①.

Update rule of Q-learning for updating the Q-function after a transition from a state i state to a state j using action u and observing immediate reward $r(i,u)$:

$$Q(S,A) \leftarrow Q(S,A) + \alpha (R + \gamma \max_{A'} Q(S',A') - Q(S,A))$$

Diagram illustrating the components of the Q-learning update rule:

- $Q(S,A)$ (left): Q-function \rightarrow new value
- $Q(S,A)$ (middle): Q-function \rightarrow old value
- α : learning factor \rightarrow "How much does new estimate effect $Q(S,A)$ "
- R : Reward
- γ : discounting factor
- $\max_{A'} Q(S',A')$: maximum Q given action
- $Q(S,A)$ (right): Q-function \rightarrow old value

How would you handle transitions to or within the goal state?

If the target is the goal state just use the reward R as value of the Q-function. Different from the lecture we would use a slightly different update rule: $Q(S,A) \leftarrow (1-\alpha) Q(S,A) + \alpha R$
 \Rightarrow Then if the target is reached the Q-function value remains the reward R of the goal state

2

8 states with 5 actions each.

Q-function → action

Q-function

0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

$$Q(S,A) \leftarrow Q(S,A) + \alpha (R + \gamma \max_{A'} Q(S',A') - Q(S,A))$$

State 1 Action 1
("down")

Action: 1 = down

2 = up

3 = left

4 = right

5 = stay

-1	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

1		0
2	3	4
5	6	7

State 2 Action 4
("right")

-1	0	0	0	0
0	0	0	-1	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

1		0
2	3	4
5	6	7

State 3 Action 2
("up")

-1	0	0	0	0
0	0	0	-1	0
0	-1	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

State 3 Action 4

-1	0	0	0	0
0	0	0	-1	0
0	-1	0	-1	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

State 4 Action 2

-1	0	0	0	0
0	0	0	-1	0
0	-1	0	-1	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

1		G
2	3	4
5	6	7

1		G
2	3	4
5	6	7

1		4
2	3	4
5	6	7

improved Q function after initial episode

