

Reinforcement Learning

an introduction

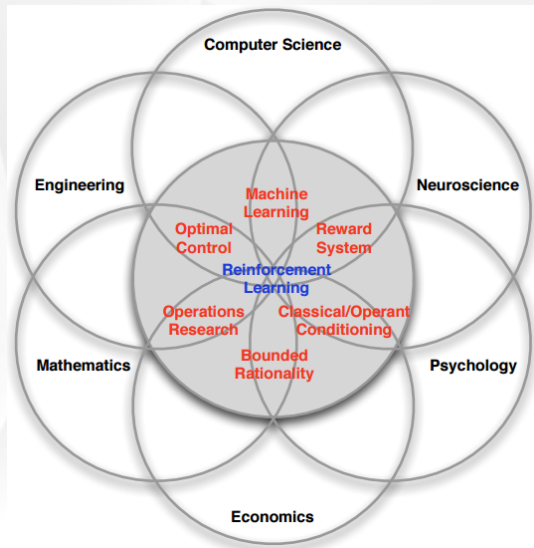
Sindre B. Remman & Thomas N. Larsen

NTNU

1 What is Reinforcement Learning?

2 Q-Learning in the Gridworld environment

3 Open problems





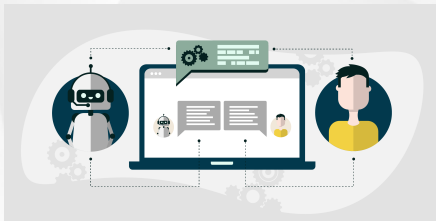
Use cases:

- AlphaGo



Use cases:

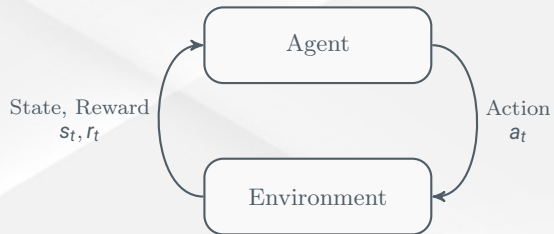
- AlphaGo
- Autonomous Cars



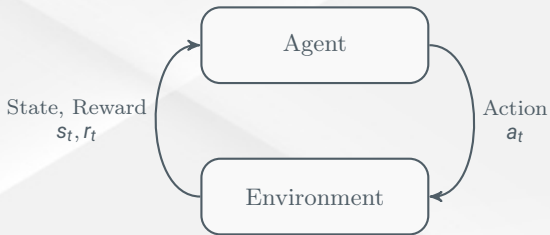
Use cases:

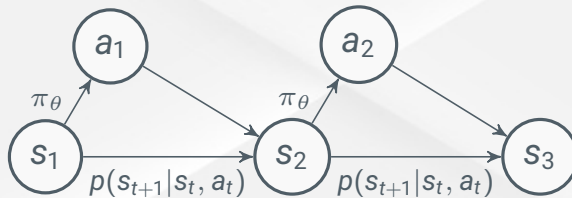
- AlphaGo
- Autonomous Cars
- Natural Language Processing

- Agent
- Environment



- Agent
 - Learner and decision-maker.
- Environment
 - Everything outside the agent that the agent interacts with





Markov Decision Process

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

- Policy Gradient methods

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

- Policy Gradient methods
- Actor-Critic methods

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

- Policy Gradient methods
- Actor-Critic methods
- Value-based methods

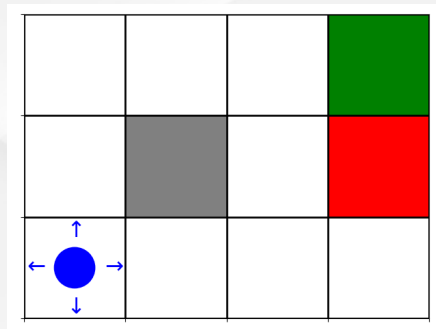
$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

- Policy Gradient methods
- Actor-Critic methods
- Value-based methods
- (Model-based methods)

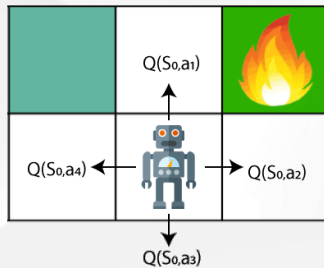
$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

- Policy Gradient methods
- Actor-Critic methods
- Value-based methods
- (Model-based methods)
- (On-policy and off-policy)

- Discrete state and action space.
 - States: 2D coordinates (x,y)
 - Actions: move (left, right, up, down)
- Markov property.
- Stochastic transition probability.
- An episode is finished when a terminal state is reached.
- Rewards:
 - Green: "Good" terminal state (reward = +1)
 - Red: "Bad" terminal state (reward = -1)
 - White: reward = 0
 - Gray: Inaccessible

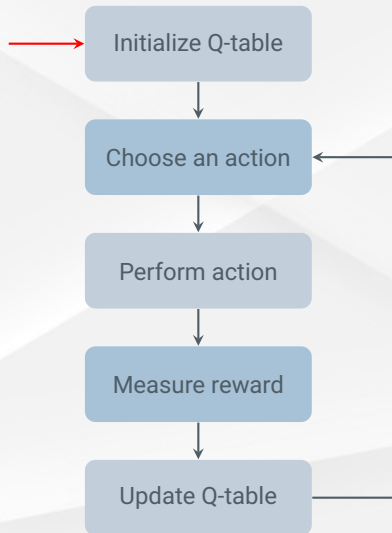


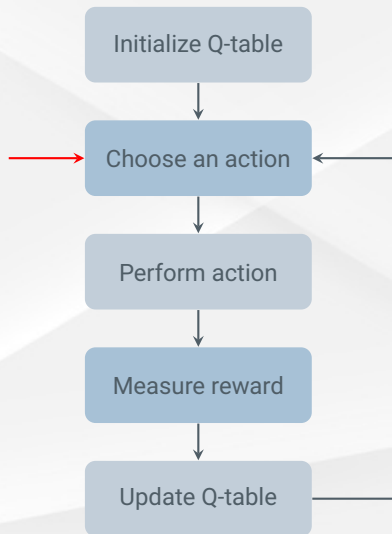
The Gridworld environment

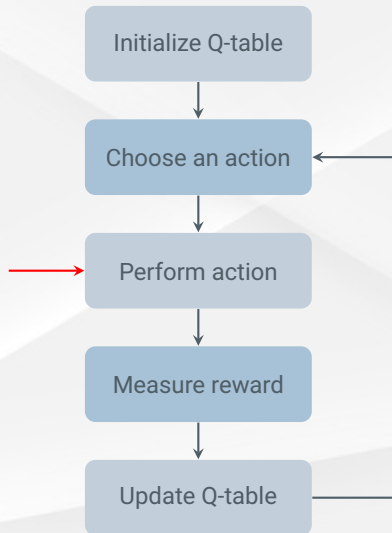


		Actions			
		A_1	A_2	...	A_M
States	S_1	$Q(S_1, A_1)$	$Q(S_1, A_2)$...	$Q(S_1, A_M)$
	S_2	$Q(S_2, A_1)$	$Q(S_2, A_2)$...	$Q(S_2, A_M)$

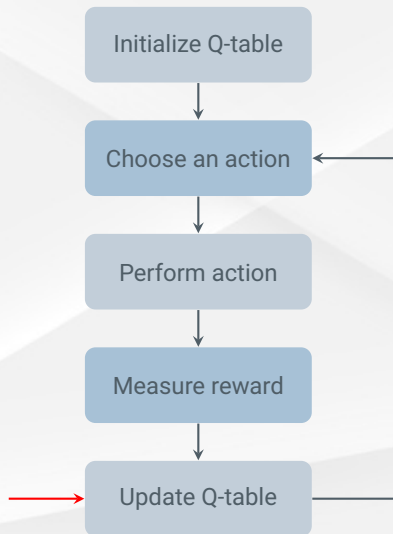
	S_N	$Q(S_N, A_1)$	$Q(S_N, A_2)$...	$Q(S_N, A_M)$

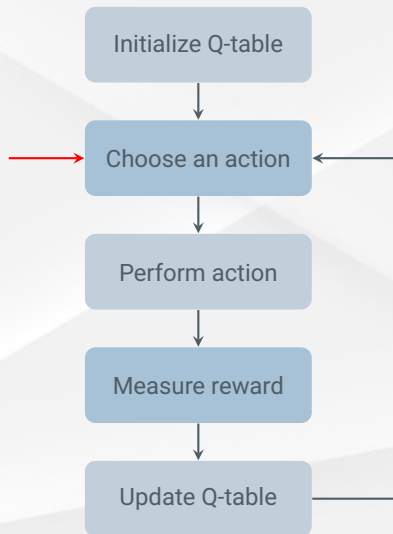


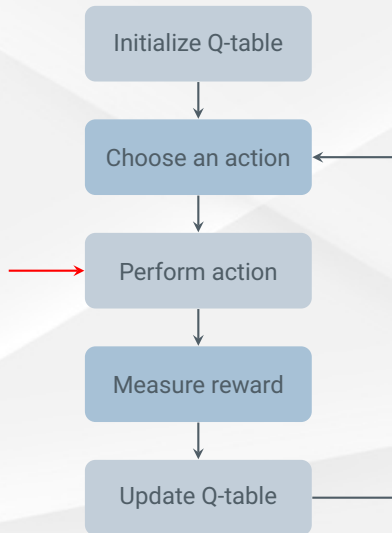


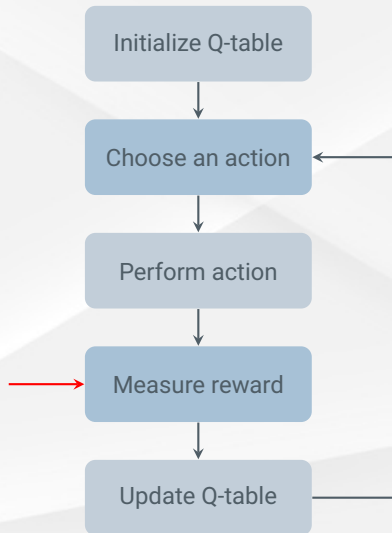


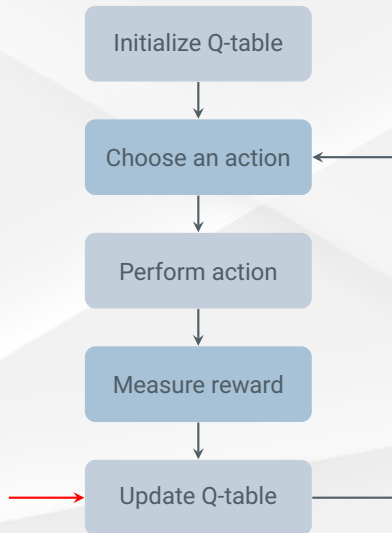












Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

0.55	0.42	0.96	
0.60 0.00	0.44 0.00	0.79 0.53	
0.72	0.65	0.38	
0.02		0.98	
0.78 0.87		0.46 0.78	
0.83		0.80	
0.52	0.46	0.61	
0.26 0.77	0.02 0.62	0.94 0.68	
0.41	0.57	0.62	0.44

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
for episode = 1, M **do**

end

0.55	0.42	0.96		
0.60 0.00	0.44 0.00	0.79 0.53		
0.72	0.65	0.38		
0.02		0.98		
0.78 0.87		0.46 0.78		
0.83		0.80		
0.52	0.46	0.61	0.36	
0.26 0.77	0.02 0.62	0.94 0.68	0.70 0.06	
0.41	0.57	0.62	0.44	

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize \mathcal{S}

end

0.0	0.0	0.0	1.0
0.0		0.0	-1.0
0.0	0.0	0.0	0.0

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize S

while $S \neq \text{terminal}$ **do**

end

end

0.0	0.0	0.0	1.0
0.0		0.0	-1.0
0.0	0.0	0.0	0.0

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

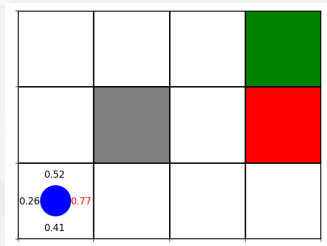
 Initialize S

while $S \neq \text{terminal}$ **do**

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

end

end



Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize S

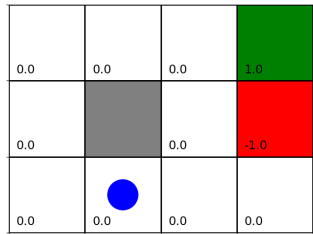
while $S \neq \text{terminal}$ **do**

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

end

end



0.0	0.0	0.0	1.0
0.0		0.0	-1.0
0.0	0.0	0.0	0.0

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize S

while $S \neq \text{terminal}$ **do**

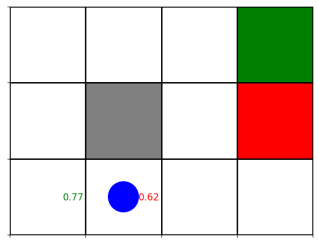
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

end

end



Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize S

while $S \neq \text{terminal}$ **do**

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

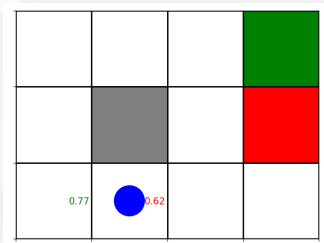
 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

end

end

$Q(\text{bottom_left}, \text{RIGHT}) \leftarrow 0.77 + 0.005 \times [0 + 0.99 \times 0.62 - 0.77]$



Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
for episode = 1, M **do**
 Initialize S
 while $S \neq \text{terminal}$ **do**
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$
 end
end

1.41	1.75	1.89	
1.45 1.45	1.63 1.89	1.86 1.91	
1.68	1.80	1.78	
1.61		1.75	
1.64 1.63		1.72 0.16	
1.69		1.67	
1.67	1.72	1.74	0.17
1.68 1.70	1.70 1.72	1.72 1.63	1.64 1.33
1.68	1.71	1.72	1.46

After 100k episodes

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
for episode = 1, M **do**
 Initialize S
 while $S \neq \text{terminal}$ **do**
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$
 end
end

1.69	1.86	1.89	
1.70 1.86	1.78 1.88	1.86 1.90	
1.73	1.86	1.77	
1.84		1.77	
1.70 1.71		1.74 0.15	
1.70		1.62	
1.71	1.70	1.74	0.14
1.69 1.72	1.70 1.73	1.71 1.64	1.63 1.50
1.69	1.70	1.72	1.61

After 200k episodes

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
for episode = 1, M **do**
 Initialize \mathcal{S}
 while $\mathcal{S} \neq \text{terminal}$ **do**
 Choose A from \mathcal{S} using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, \mathcal{S}'
 $Q(\mathcal{S}, A) \leftarrow Q(\mathcal{S}, A) + \alpha[R + \gamma \max_a Q(\mathcal{S}', a) - Q(\mathcal{S}, A)]$
 $\mathcal{S} \leftarrow \mathcal{S}'$
 end
end

1.85	1.87	1.89	
1.84 1.86	1.85 1.89	1.86 1.91	
1.83	1.87	1.77	
1.84		1.78	
1.82 1.82		1.71 0.09	
1.81		1.69	
1.82	1.78	1.75	0.11
1.80 1.78	1.80 1.76	1.78 1.65	1.66 1.55
1.80	1.78	1.75	1.60

After 300k episodes

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$
 Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$
for episode = 1, M **do**
 Initialize \mathcal{S}
 while $S \neq \text{terminal}$ **do**
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$
 end
end

1.85	1.87	1.89	
1.84 1.86	1.85 1.89	1.87 1.91	
1.83	1.87	1.78	
1.84		1.77	
1.83 1.83		1.75 0.13	
1.81		1.68	
1.82	1.78	1.76	0.10
1.81 1.79	1.80 1.77	1.78 1.66	1.66 1.51
1.80	1.78	1.76	1.61

After 400k episodes

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, $\gamma \in [0, 1]$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

for episode = 1, M **do**

 Initialize S

while $S \neq \text{terminal}$ **do**

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

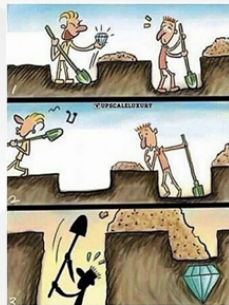
$S \leftarrow S'$

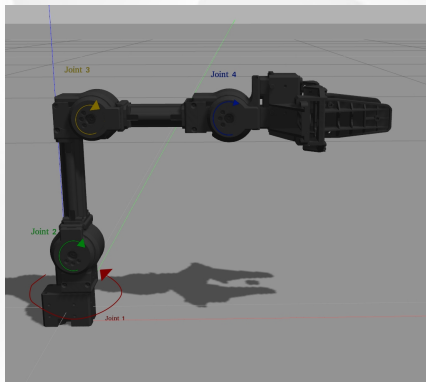
end

end

→	→	→	
↑		↑	
↑	←	←	←

Final policy







Thank you for your attention!