Introduction to Machine Learning

Lecture 16: Elementary Reinforcement Learning – Stochastic Environment Dec 9, 2019

Jie Wang

Machine Intelligence Research and Applications Lab Department of Electronic Engineering and Information Science (EEIS)

http://staff.ustc.edu.cn/~jwangx/

jiewangx@ustc.edu.cn





Stochastic Environment

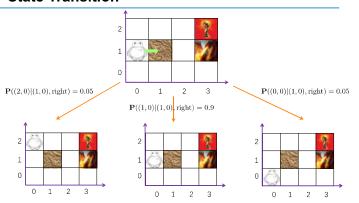
State Transition



State transition probabilities:

After the agent picks and performs a certain action, there are four possibilities for the next state: the destination state, the current state, the states to the right and left of the current state. If the states are reachable, the corresponding probabilities are 0.8, 0.1, 0.05, and 0.05, respectively; otherwise, the agent stays where it is.

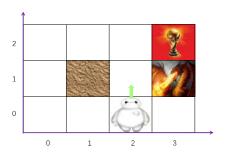
State Transition



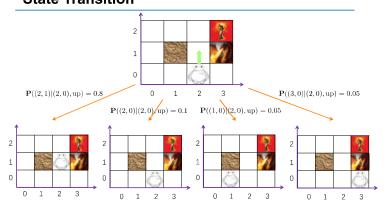
Contents

- Stochastic Environment
- Planning Algorithms
- Learning Algorithms

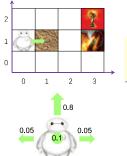
Grid World



State Transition



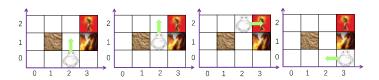
Reward



Reward:

After the agent picks and performs a certain action at its current state, it receives rewards of 100, -100, and 0, if it arrives at states (3,2), (2,2), and all the other states, respectively.

Reward



 $\mathbf{E}[r((2,0),\mathrm{up})] = 0.8 \times 0 + 0.1 \times 0 + 0.05 \times 0 + 0.05 \times 0 = 0$

 $\mathbf{E}[r((2,1),\mathrm{up})] = 0.8 \times 0 + 0.1 \times 0 + 0.05 \times -100 + 0.05 \times 0 = -5$

 $\mathbf{E}[r((2,2), \text{right})] = 0.8 \times 100 + 0.1 \times 0 + 0.05 \times 0 + 0.05 \times 0 = 80$

 $\mathbf{E}[r((3,0), \text{left})] = 0.8 \times 0 + 0.1 \times 0 + 0.05 \times -100 + 0.05 \times 0 = -5$

Value Function

- Suppose that a policy π is given.
- Starting from an arbitrary state s_t , the expected cumulative reward by following π is

$$V^{\pi}(s_t) := \mathbf{E}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s_t] = \mathbf{E}\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i} | S_t = s_t\right]$$

$$s_{t} \xrightarrow{a_{t}} s_{t+1} \xrightarrow{a_{t+1}} s_{t+1} \xrightarrow{s_{t+2}} a_{t+2} \xrightarrow{r_{t+2}} \cdots$$

$$a_{t} = \pi(s_{t}) \quad r_{t} = r(s_{t}, a_{t}) \quad s_{t+1} = \delta(s_{t}, a_{t})$$
random variable

Value Function

Tower property

$$\mathbf{E}[X|Y] = \mathbf{E}[\mathbf{E}[X|Y,Z]|Y]$$

· A simpler version

$$\mathbf{E}[X] = \mathbf{E}[\mathbf{E}[X|Z]]$$

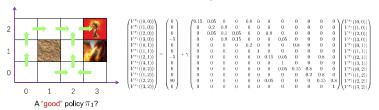
• Example: how to find the average height of Chinese men?

$$E[\text{height}] = E\big[E[\text{height}|\text{province}]\big] = \sum_{\text{province}} P(\text{province})E[\text{height}|\text{province}]$$

Value Function - Bellman Equation

Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s,\pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s,\pi(s)) V^{\pi}(s')$$



Markov Decision Process (MDP)

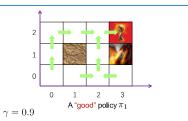
- Indeed, we have already introduced the so-called MDP, which is defined (rigorously) by
 - \bullet a set of states ${\cal S}$, possibly infinite

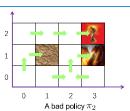
a set of states \mathcal{A} , possibly infinite a set of actions \mathcal{A} , possibly infinite

MRT Chapter 14

- an initial state $s_0 \in \mathcal{S}$
- a transition probability $\mathbf{P}[s'|s,a]$: distribution over destination states $s'=\delta(s,a)$
- ullet a reward probability $\mathbf{P}[r|s,a]$: distribution over rewards r'=r(s,a)
- This model is Markovian because the transition and reward probabilities only depend on the current state and the action picked and performed at the current state, instead of the previous sequence of states and actions performed.
- · In this lecture, we assume that
 - . the states and the actions are finite

Value Function





How to find V^{π_1} and V^{π_2} ?

Value Function - Bellman Equation

• Starting from an arbitrary state s_t , the expected cumulative reward by following π is

$$V^{\pi}(s_t) := \mathbf{E}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s_t] = \mathbf{E}\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i} | S_t = s_t\right]$$

$$s_t \xrightarrow{a_t} s_{t+1} \xrightarrow{a_{t+1}} s_{t+1} \xrightarrow{s_{t+2}} s_{t+2} \xrightarrow{r_{t+2}} \dots$$

$$a_t = \pi(s_t) \quad r_t = r(s_t, a_t) \quad s_{t+1} = \delta(s_t, a_t)$$
random variable

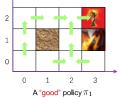
Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s, \pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s, \pi(s))V^{\pi}(s')$$

Value Function - Bellman Equation

Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s, \pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s, \pi(s)) V^{\pi}(s')$$



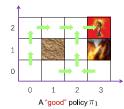


 $V = R + \gamma TV$

Value Function - Bellman Equation

Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s,\pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s,\pi(s)) V^{\pi}(s')$$



$$V = R + \gamma TV$$

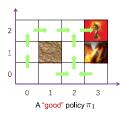
$$\downarrow$$

$$V = (I - \gamma T)^{-1}.$$

Value Function - Bellman Equation

Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s,\pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s,\pi(s)) V^{\pi}(s')$$



Theorem: For a finite MDP, Bellman's equation admits a unique solution that is given by

$$V = (I - \gamma T)^{-1}R$$

- The $\operatorname{vector} R$ and $\operatorname{matrix} T$ depend on the policy

The Q Function

- Learning the optimal policy is challenging
- An alternative approach to find the optimal policy indirectly is by computing the state-action value function (Q function)

$$Q(s,a) = \mathbf{E}[r(s,a)] + \gamma \sum_{s} \mathbf{P}(s'|s,a) V^*(s')$$

Q(s,a) is the expected accumulated reward by performing

· The definition of the optimal policy implies that

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$$

Notice that

$$V^*(s) = \max_{a} Q(s, a)$$

· All together, we have

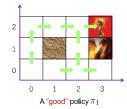
$$Q(s,a) = \mathbf{E}[r(s,a)] + \gamma \sum_{s'} \left[\mathbf{P}(s'|s,a) \max_{a'} Q(s',a') \right]$$
 Bellman Equation

Planning Algorithms

Value Function - Bellman Equation

Bellman Equation

$$V^{\pi}(s) = \mathbf{E}[r(s,\pi(s))] + \gamma \sum_{s'} \mathbf{P}(s'|s,\pi(s)) V^{\pi}(s')$$

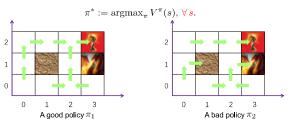


$$V=R+\gamma TV$$

$$V=\underbrace{\left(I-\gamma T\right)}_{\text{invertible ?}}^{-1}R$$

The Learning Task Revisited

• The learning task for RL scenarios is to learn an optimal policy in the sense that



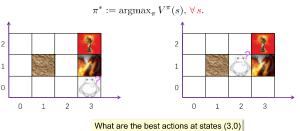
• For π_1 and π_2 , we have

$$V^{\pi_1}(s) \ge V^{\pi_2}(s), \,\forall \, s.$$

Indeed, \(\pi_1\) is the optimal policy.

Quiz

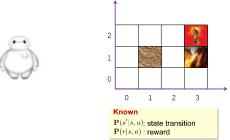
• The learning task for RL scenarios is to learn an optimal policy in the sense that



and (2,1), i.e., $\pi^*((3,0))$ and $\pi^*((2,1))$?

Planning

 Planning: to find the optimal policy, there is no need for the agent to actually perform actions and interact with the environment



Value Iteration

· Value iteration aims to find the optimal value function and thus the optimal policy

```
\begin{split} \textbf{Initialize} \ & V(s) \ \text{to arbitrary values} \\ \textbf{while} \ & \textbf{termination conditions does not hold} \\ \textbf{For} \ & s \in \mathcal{S} \\ \textbf{For} \ & a \in \mathcal{A} \\ & Q(s,a) \leftarrow \mathbf{E}[r(s,a)] + \gamma \sum_{s'} \mathbf{P}(s'|s,a) V(s') \\ & V(s) \leftarrow \max_{a} \, Q(s,a) \end{split}
```

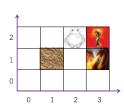
Value Iteration

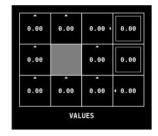
· Value iteration aims to find the optimal value function and thus the optimal policy



Value Iteration

Value iteration aims to find the optimal value function and thus the optimal policy





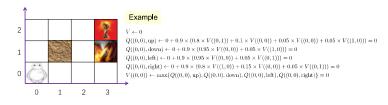
Policy Iteration

· Policy iteration improves the policy directly

```
\begin{split} \textbf{Initialize} & \ \pi \leftarrow \pi_2, \pi' \neq \pi_2 \\ \textbf{while}(\pi \neq \pi') & \ V \leftarrow (I - \gamma T^\pi)^{-1} R^\pi \\ & \ \pi' \leftarrow \pi \\ & \ \textbf{For} \ s \in \mathcal{S} \\ & \ \pi(s) \leftarrow \text{argmax}_a \ \textbf{E}[r(s,a)] + \gamma \sum_{s'} \textbf{P}(s'|s,a) V(s') \end{split}
```

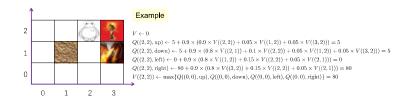
Value Iteration

· Value iteration aims to find the optimal value function and thus the optimal policy



Value Iteration

· Value iteration aims to find the optimal value function and thus the optimal policy



Value Iteration

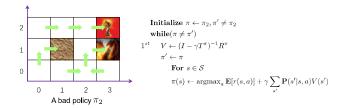
· Value iteration aims to find the optimal value function and thus the optimal policy

Theorem: For any initial value V, the sequence generated by the value iteration algorithm converges to V^* .

The key to the proof is the contraction mapping theorem

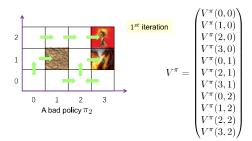
Policy Iteration

· Policy iteration improves the policy directly



Policy Iteration

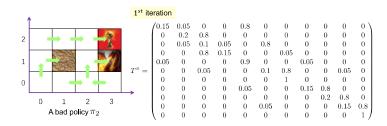
· Policy iteration improves the policy directly



11 states in total

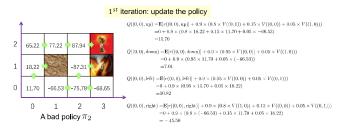
Policy Iteration

· Policy iteration improves the policy directly



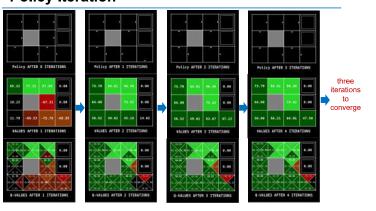
Policy Iteration

Policy iteration improves the policy directly



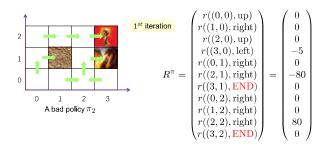
 $\pi((0,0)) = \operatorname{argmax}_{\{\text{up, down, left, right}\}} \{Q((0,0), \operatorname{up}), Q((0,0), \operatorname{down}), Q((0,0), \operatorname{left}), Q((0,0), \operatorname{right})\}$

Policy Iteration



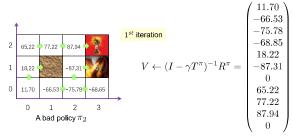
Policy Iteration

· Policy iteration improves the policy directly



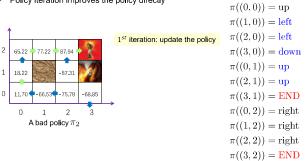
Policy Iteration

· Policy iteration improves the policy directly

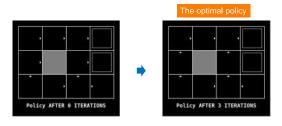


Policy Iteration

· Policy iteration improves the policy directly



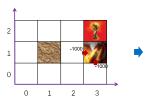
Policy Iteration

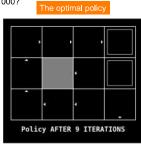


Is this an always winning policy?

Policy Iteration

• What if the reward for getting into (3,1) is -1000?



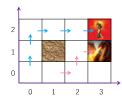


This is an always winning policy (why?).

Learning Algorithms

Learning

- Learning: as the environment model, i.e., the transition and reward probabilities, is unknown, the agent may need to learn them based on the training information.
 - · Model-free approach: the agent learns the optimal policy directly, e.g., Q-learning
 - Model-based approach: the agent first learns the environment model and then the
 optimal policy



Examples of training data

$$(0,0) \frac{up}{0} (0,1) \frac{up}{0} (0,2) \frac{right}{0} (1,2) \frac{right}{0} (2,3) \frac{right}{100} (3,2)$$

$$(1,0) \frac{right}{0} (2,0) \frac{up}{0} (2,1) \frac{right}{-100} (3,1)$$

The Q-learning Algorithm

Recall the Q-learning algorithm for the deterministic environment

- Initialize the matrix \hat{Q} to zero
- Observe the current state s
- Do forever:
 - Pick and perform an action \boldsymbol{a}
 - Receive immediate reward r
 - Observe the new state s'
 - Update

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

s ← s'

A sufficient condition for $\hat{Q}(s,a)$ to converge is to visit each state-action pair infinitely often

Policy Iteration



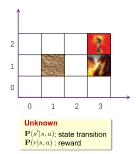


- · For the same task, different reward strategies lead to different optimal policy
- · According to your preference, you need to carefully design your reward strategy

Learning

 Learning: as the environment model, i.e., the transition and reward, is unknown, the agent may need to learn them based on the training information.





Nondeterministic Rewards and Actions





How to find the optimal policy without the state transition and reward probabilities?

The Q-learning Algorithm

 For the stochastic environment, we replace the random variables with their expectations in the definition of Q values.



Q: How to find the expectation of a random variable *X*?

A: Keep sampling and recording its running average



The Q-learning Algorithm

Initialize \hat{Q} arbitrarily

For all episodes

Initialize s

Alpaydin 2014, Chapter 18

Repeat

Choose a using policy derived from Q, e.g., ϵ -greedy

Take action a, observe r and s'

Update $\hat{Q}(s, a)$:

$$\alpha_n = \frac{1}{1 + n((s,a))}$$
 the number of visits of (s,

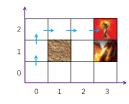
$$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha_n(r + \gamma \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a))$$

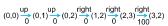
 $s \leftarrow s'$

Until s is goal state

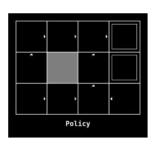
There are other ways to select α_n to guarantee that \hat{Q} converges to its optimal value. Mitchell 1997, Chapter 13

The Q-learning Algorithm



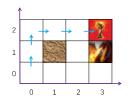


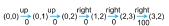
- an example epicode
- an example episode
 the initial state in each episode could NOT be fixed (why?)



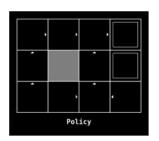
 $\epsilon = 0.3$

The Q-learning Algorithm



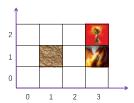


- · an example episode
- the initial state in each episode could NOT be fixed (why?)



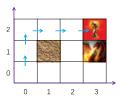
 $\epsilon = 0.3$

SARSA





The Q-learning Algorithm



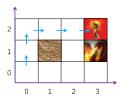
$$(0,0) \xrightarrow{\mathsf{up}} (0,1) \xrightarrow{\mathsf{up}} (0,2) \xrightarrow{\mathsf{right}} (1,2) \xrightarrow{\mathsf{right}} (2,3) \xrightarrow{\mathsf{right}} (3,2)$$

- an example episode
- the initial state in each episode could NOT be fixed (why?)



 $\epsilon = 0.3$

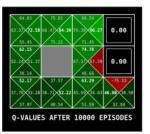
The Q-learning Algorithm



$$(0,0) \xrightarrow{up} (0,1) \xrightarrow{up} (0,2) \xrightarrow{right} (1,2) \xrightarrow{right} (2,3) \xrightarrow{right} (3,2)$$

an example episodethe initial state in each episode

could NOT be fixed (why?)



$$\epsilon = 0.3$$

SARSA

Initialize $\hat{Q} \leftarrow 0$

For all episodes

 ${\bf Initialize}\ s$

Choose a using policy derived from Q, e.g., ϵ -greedy

 ${\bf Repeat}$

Take action a, observe r and s'

Choose a' using policy derived from Q, e.g., ϵ -greedy

Update $\hat{Q}(s, a)$:

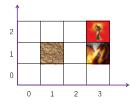
$$\alpha_n = \frac{1}{1 + n((s, a))}$$

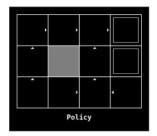
$$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha_n(r + \gamma \hat{Q}(s',a') - \hat{Q}(s,a))$$

 $s \leftarrow s', a \leftarrow a'$ Until s is goal state

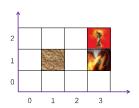
Alpaydin 2014, Chapter 18

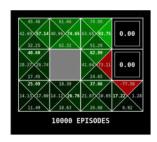
SARSA

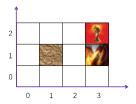


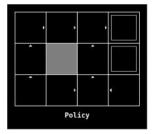


SARSA SARSA









Questions

