

Reinforcement Learning

基本名词

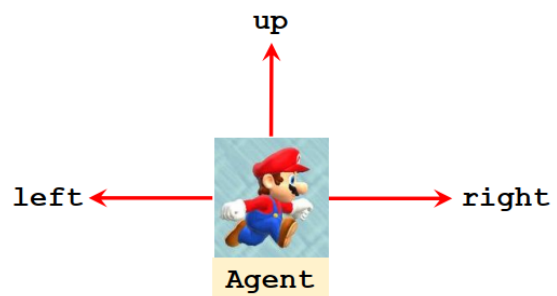
State and Action

Terminology: **state** and **action**

state s (this frame)



Action $a \in \{\text{left, right, up}\}$



state是场景

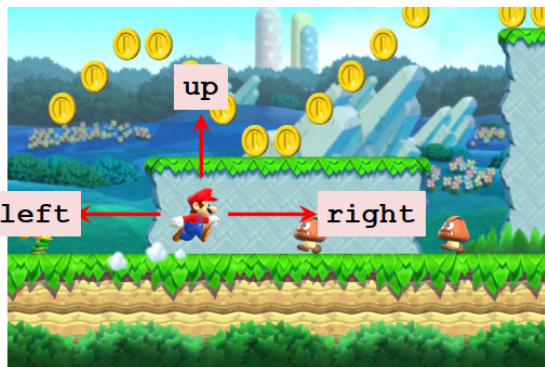
agent为智能体，也就是执行操作的对象，此处是马里奥，也可以是机器人，车等

action为agent做出的动作

Policy

Terminology: **policy**

policy π



- $\pi(a | s)$ is the probability of taking action $A = a$ given state s , e.g.,
 - $\pi(\text{left} | s) = 0.2$,
 - $\pi(\text{right} | s) = 0.1$,
 - $\pi(\text{up} | s) = 0.7$.
- Upon observing state $S = s$, the agent's action A can be random.

Policy函数，即策略函数，通常是一个概率分布函数，如该ppt中，在当前state下，马里奥选择往左走的策略的概率为0.2。该state下马里奥有三种策略，会在其中随机抽样，随机性使得policy更加灵活，难以被预测

Reward

Terminology: reward

reward R



- Collect a coin: $R = +1$
- Win the game: $R = +10000$
- Touch a Goomba: $R = -10000$ (game over).
- Nothing happens: $R = 0$

reward为奖励函数，如此处吃到金币+1分，通关+1w分，碰到敌人（game over）扣1w分

State Transition

Terminology: state transition

state transition



- E.g., “up” action leads to a new state.
- State transition can be random.
- Randomness is from the environment.
- $p(s'|s, a) = \mathbb{P}(S' = s' | S = s, A = a)$.

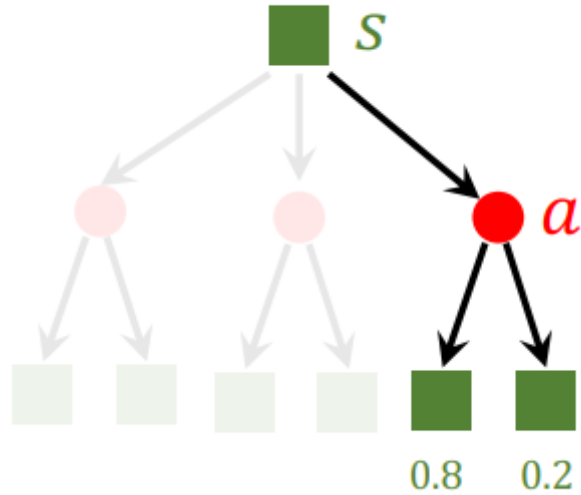
状态转移，顾名思义是old state变为new state，上图P为条件概率密度函数，意思是如果观测到当前的状态S以及动作A，下一个状态变成S一撇的概率

由于env以及action的随机性，所以状态转移具有随机性

Two Sources of Randomness

整个过程中的随机性主要来源于两点：

1. 是action的随机性，这个很好理解
2. 是state的随机性，在马里奥的例子中可以理解为敌人移动也是随机的，无法被agent知晓



- The randomness in **action** is from the policy function:

$$A \sim \pi(\cdot | s).$$

- The randomness in **state** is from the state-transition function:

$$S' \sim p(\cdot | s, a).$$

Trajectory

整个过程可以认为是

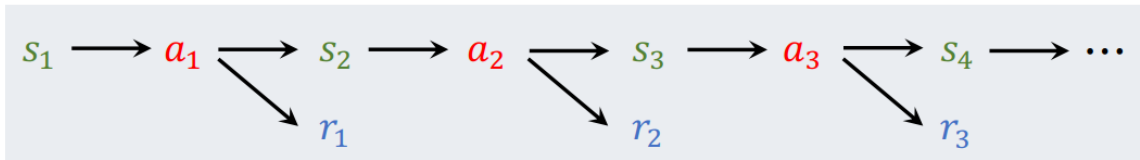
1. 观察state
2. 做出action
3. 观察到新的state并得到奖励（或惩罚）

Play game using AI

- (state, action, reward) trajectory:

$$\underline{s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.}$$

- One episode is from the beginning to the end (Mario wins or dies).



Rewards and Returns

returns 定义为未来所有 cumulative future reward 未来的累积奖励

$$U_t = R_t + R_{t+1} + R_{t+2} + R_{t+3} + \dots$$

但其实之后的奖励对当前时刻不是同等重要的，如选择现在给你100元和一年后给你100元，大部分人会选择立刻得到一百

所以引入折扣回报 Discounted return

Definition: Discounted return (at time t).

$$U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{n-t} R_n.$$

gamma 介于0到1，为超参数

R 与 S 和 A 有关

- Reward R_i depends on S_i and A_i .
- States can be random: $S_i \sim p(\cdot | s_{i-1}, a_{i-1})$.
- Actions can be random: $A_i \sim \pi(\cdot | s_i)$.
- If either S_i or A_i is random, then R_i is random.

所以U 与未来所有 S A 有关

At time t , the rewards, R_t, \dots, R_n , are random, so the return U_t is random.

- Reward R_i depends on S_i and A_i .
- U_t depends on R_t, R_{t+1}, \dots, R_n .
- $\rightarrow U_t$ depends on $S_t, A_t, S_{t+1}, A_{t+1}, \dots, S_n, A_n$.

Value Function

价值函数

Action-value function

U_t 其实是一个随机变量，他依赖于之后的所有动作和状态， t 时刻我们并不知道 U_t 是什么，所以我们可以对 U_t 求期望得到一个函数，即Action-value function，动作价值函数， Q_{π}

Definition: Action-value function for policy π .

$$Q_{\pi}(s_t, a_t) = \mathbb{E}[U_t | S_t = s_t, A_t = a_t].$$

- Return U_t depends on actions $A_t, A_{t+1}, A_{t+2}, \dots$ and states $S_t, S_{t+1}, S_{t+2}, \dots$
- Actions are random: $\mathbb{P}[A = a | S = s] = \pi(a|s)$. (Policy function.)
- States are random: $\mathbb{P}[S' = s' | S = s, A = a] = p(s'|s, a)$. (State transition.)

此时我们 U_t 未知， S_t 和 A_t 是变量，且他们的概率分布已知（ S_t, A_t 分布），则可以用积分的方式把将来的SA对当前时刻 U_t 的影响通过积分求期望的方式获得

将 t 时刻之后的随机变量 A, S 都用积分积掉，之后得到的 Q_{π} 就只与当前时刻 t 的SA以及 π 有关

Action-value function的实际意义，一直policy函数 π 以及当前 t 时刻的 s ，则可以通过Action-value function Q_{π} 对每个action打分，看做出哪个action，最终 U_t 的期望最高

Optimal Action-value function

最优动作价值函数

之前说的Action-value function与 π, SA 有关，而我们有很多个policy函数 π ，我们要使得Action-value function最大，则可以做一个动态规划，取得最优的policy函数 π ，使得Action-value function最大，这样就可以消除 π 对Action-value function的影响（因为policy函数 π 已经确定了，不再是变量），得到一个最优的Action-value function，即Optimal Action-value function

Definition: Optimal action-value function.

$$Q^*(s_t, a_t) = \max_{\pi} Q_{\pi}(s_t, a_t).$$

Definition: State-value function.

$$\bullet V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)]$$

Definition: State-value function.

- $V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)] = \sum_a \pi(a|s_t) \cdot Q_{\pi}(s_t, a)$. (Actions are discrete.)
- $V_{\pi}(s_t) = \mathbb{E}_A [Q_{\pi}(s_t, A)] = \int \pi(a|s_t) \cdot Q_{\pi}(s_t, a) da$. (Actions are continuous.)

将 Q_{π} 对 A 求期望（将 A 当作随机变量），从而消掉 A

物理意义在于可以评估目前状态的胜算

Conclusion

Understanding the Value Functions

- **Action-value function:** $Q_{\pi}(s, a) = \mathbb{E} [U_t | S_t = s, A_t = a]$.
- Given policy π , $Q_{\pi}(s, a)$ evaluates how good it is for an agent to pick action a while being in state s .
- **State-value function:** $V_{\pi}(s) = \mathbb{E}_A [Q_{\pi}(s, A)]$
- For fixed policy π , $V_{\pi}(s)$ evaluates how good the situation is in state s .
- $\mathbb{E}_S [V_{\pi}(S)]$ evaluates how good the policy π is.

How does AI control the agent

1. policy-based learning 策略学习：已知 π ， S ，可以求得每个 A 的概率，再随机抽样
2. value-based learning 价值学习：求Optimal Action-value function Q-star

How does AI control the agent?

Suppose we have a good policy $\pi(a|s)$.

- Upon observe the state s_t ,
- random sampling: $a_t \sim \pi(\cdot | s_t)$.

Suppose we know the optimal action-value function $Q^*(\underline{s}, \underline{a})$.

- Upon observe the state s_t ,
- choose the **action** that maximizes the value: $a_t = \operatorname{argmax}_a Q^*(s_t, a)$.