# ch18: Reinforcement Learning

## Introduce to RL

### reinforcement learning 和其他的差异（P4）

| Supervised Learning | Unsupervised Learning | Reinforcement Learning |
|---|---|---|
| **Data**: $(x, y)$ $x$ is data, $y$ is label | 只考虑到可以学 **Data**: $x$ 习到的规律 $x$ is data, no label | 只考虑数据带来 **Data**: *state-action* 的收益 pairs + rewards |
| **Goal**: Learn to map $x \rightarrow y$ | **Goal**: Learn underline structure. | **Goal**: Learn from interacting with environment. |
| **Examples**: Classification, regression, etc. | **Examples**: Clustering, dim reduction, etc. | **Examples**: Robot, Game, etc. |
| This thing is an apple. | This thing is like the other thing. | Eat this thing because it will keep you alive. |

### Reinforcement Learning: Key Concepts（P9-12）

**Agent**: an actor which takes actions *and learn knowledge*
**Environment**: the world in which the agent exists and operates.
也会为agent 提供 feedback

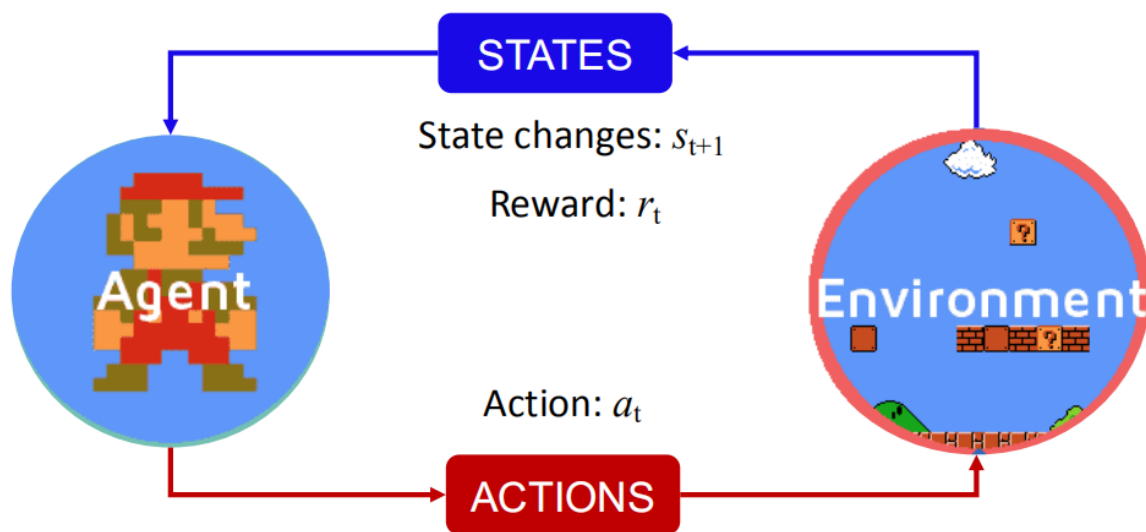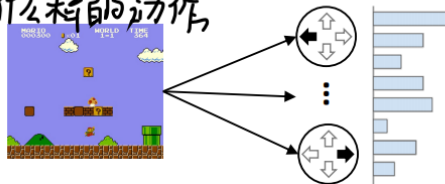**Action**: a move the agent can make in the environment.
**Action space A**: the set of possible actions an agent can make in the environment.

**State:** observations  agent 观察到环境的状态
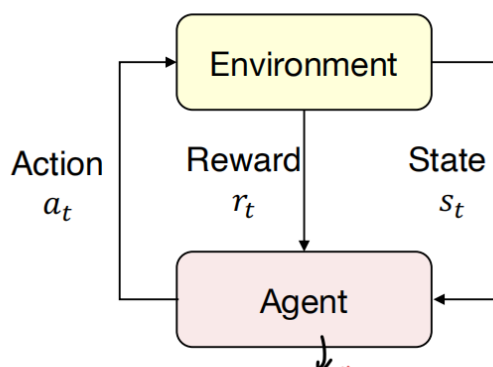**Reward:** feedback that scores the agent's action.

**Policy** $\pi$: S → A: a function that maps from state to action

在一个 state 下面应该采取什么样的动作，是一个概率分布



State changes: $s_{t+1}$

Reward: $r_t$

Action: $a_t$

## Reinforcement Learning Problem Overview（P13）

- Learning from **interacting** with an **environment**.
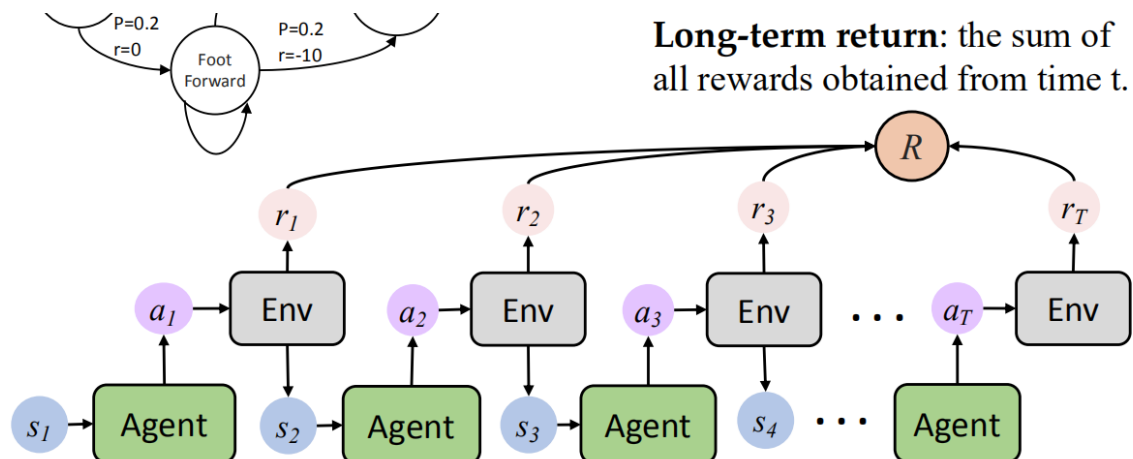


Action $a_t$

Reward $r_t$

State $s_t$

- a set of environment states S;
- a set of actions A;
- rules of transitioning between states; 获得reward的规则
- rules that determine the scalar immediate reward of a transition;
- rules that describe what the agent observes.
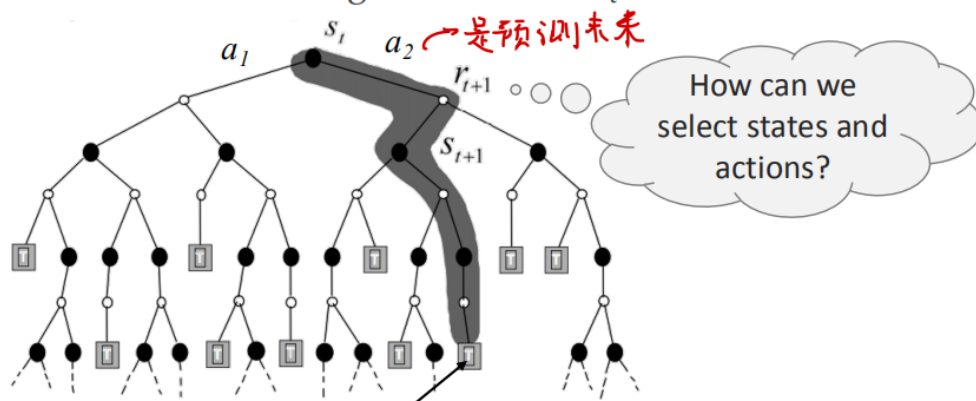
获得 reward 的时候还会发生状态的转移

目标
**Objective:** learning a policy (i.e., state-action mapping) that maximizes the long-term payoff. $\sum r_t$ 最大

# Reinforcement Learning Process（P15）



P=0.2 r=0
Foot Forward
P=0.2 r=-10

**Long-term return**: the sum of all rewards obtained from time t.

# In the View of State Space（P16）

- The agent finds a path in the state space that potentially leads to the maximum long-term return $R_t$.



是预测未来

How can we select states and actions?

Actual long-term return following $s_t$: $R_t = \sum_{i=t}^{\infty} r_i = r_t + r_{t+1} + \dots$

Discounted long-term return: $R_t = \sum_{i=t}^{\infty} \gamma^i r_i = r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} + \dots$

每个动作对于越未来的
收益越有折扣，首先考虑眼前

# Scoring (s, a) pairs: Q-function（P17）

- The **Q-function** captures the **expected total future reward** an agent in state s can receive by executing a certain action a
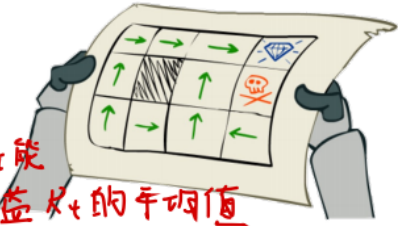
$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$$

score(state, action)

在 $s_t$ 下采取 $a_t$ 就能获到的收益 $R_t$ 的平均值

$$R_t = \sum_{i=t}^{\infty} \gamma^i r_i = r_t + \gamma r_{t+1} \ldots + \gamma^n r_{t+n} + \ldots$$

$$\begin{aligned}
Q(s_t, a_t) &= \mathbb{E}[R_t | s_t, a_t] \\
&= r_t + \gamma r_{t+1} \ldots + \gamma^n r_{t+n} + \ldots \\
&= r_t + \gamma (r_{t+1} \ldots + \gamma^{n-1} r_{t+n} + \ldots) \\
&= r_t + \gamma Q(s_{t+1}, a_{t+1})
\end{aligned}$$

Recursion

**根据Q-function确定policy(P18)**

$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$$

间接定义了策略

- Having obtained the Q-values, the agent derives the optimal **policy π(s)**, to infer the **best action to take** at state s

$$\pi^*(s) = \arg\max_a Q(s, a)$$

The policy should choose an action that maximizes future reward.

# Reinforcement Learning Algorithms

Value Learning

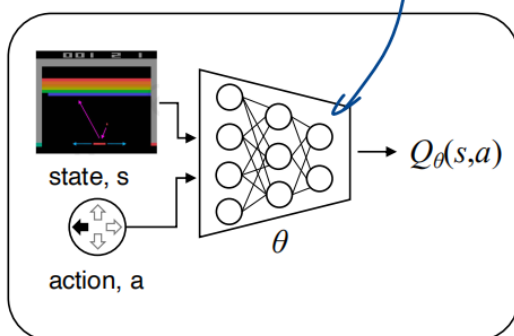Find Q(s, a)
a = argmax$_a$ Q(s, a)
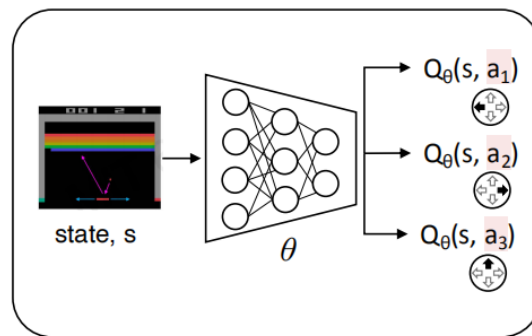
Policy Learning

Find $\pi$(s)
Sample a ~ $\pi$(s)

## PART1

## Deep-Q Networks (DQN)(P22-27)



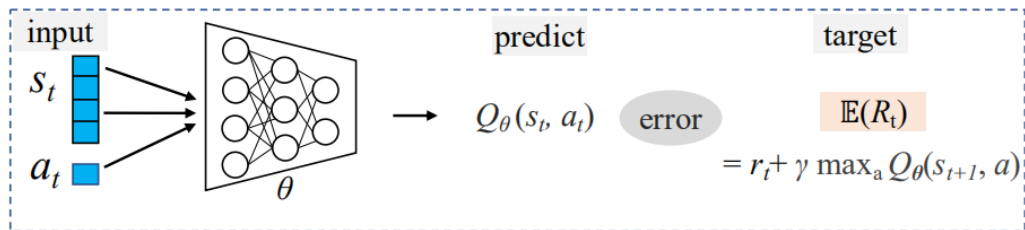- **Idea:** use deep neural networks to model the Q-function

state, s

action, a

$Q_\theta(s,a)$

$\theta$

Action + State → Expected return

state, s

$\theta$

$Q_\theta(s, a_1)$

$Q_\theta(s, a_2)$

$Q_\theta(s, a_3)$

State → Expected return for each action

What happens if we take all the best actions?
Maximize the target return → Train the agent

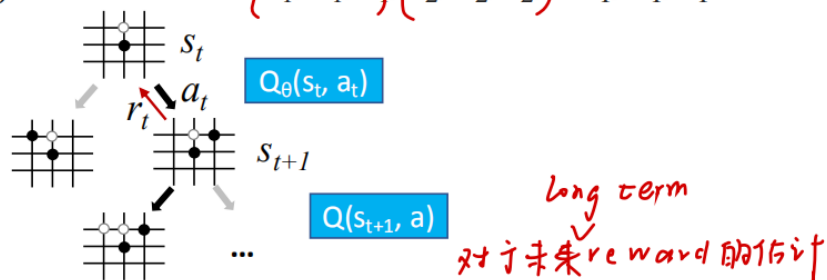Use deep neural networks (e.g., CNN) to model the Q-function.



- Train:
  - estimate the parameters $\theta$ in the $Q_\theta(s, a)$ network using agent history.
- Test:
  - calculate $Q_\theta(s, a)$ for all $a$'s under $s$ and choose $a = \text{argmax}_a Q_\theta(s, a)$.

## Training(P25-26)

- **Episode**: run agent and obtain $s_1, a_1, r_1, (s_2, a_2, r_2)..s_T, a_T, r_T$.
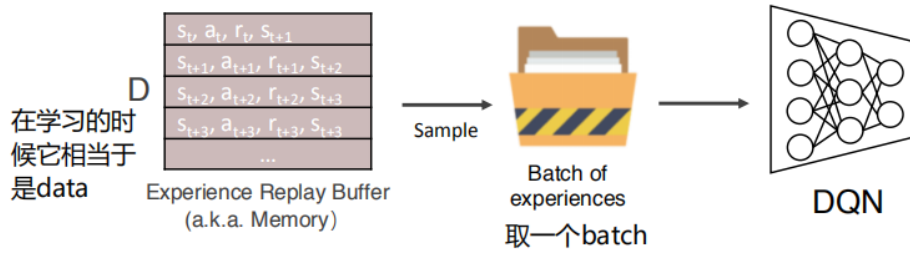
回合



long term
对于未来 reward 的估计

- For any given input $s_t$ and $a_t$, the model estimates a score $Q_\theta(s_t, a_t)$. 预测 Let $r_t$ be the reward from the environment and $s_{t+1}$ be the next state. The expected long-term reward is: $r_t + \gamma \max_{a \in A} Q(s_{t+1}, a)$.
- Our goal is to minimize their difference, i.e., the MSE loss:

$$\ell(\theta \mid s_t, a_t, r_t) = \| (r_t + \gamma \max_a Q_\theta(s_{t+1}, a) - Q_\theta(s_t, a_t) \|^2$$

target          predicted

- **Experience Replay**: run agent multiple episodes and store the transitions in a **replay memory**: $D = \{(s^{(1)}, a^{(1)}, r^{(1)}, s'^{(1)}), \ldots, (s^{(N)}, a^{(N)}, r^{(N)}, s'^{(N)})\}$



D
在学习的时候它相当于是data

| $s_t, a_t, r_t, s_{t+1}$ |
| $s_{t+1}, a_{t+1}, r_{t+1}, s_{t+2}$ |
| $s_{t+2}, a_{t+2}, r_{t+2}, s_{t+3}$ |
| $s_{t+3}, a_{t+3}, r_{t+3}, s_{t+3}$ |
| ... |

Experience Replay Buffer (a.k.a. Memory)

Sample →

Batch of experiences
取一个batch

DQN

- For multiple $(s, a, r, s') \in D$:

$$L\,(\boldsymbol{\theta}\mid D) = \mathbb{E}_{(s,a,r,s')\sim D}[\|\,(r + \gamma\,\max_{a'} Q_\theta(s',a') - Q_\theta(s,a)\,\|^2]$$

mini-batch experience replay   Maximum possible Q-value for the next state (=Q_target)   Current predicted Q-value

- Gradients:

两个 Q 都在 learning

$$\boxed{\nabla_\theta L = E_{(s,a,r,s')\sim D}\left(r + \gamma \max_{a'} Q(s',a';\theta) - Q(s,a;\theta)\right)\nabla_\theta Q(s,a;\theta)}$$

mini-batch experience replay   Q target   varying target?

- **Solution**: use some old, fixed parameters $\theta_{old}$ as a fixed Q-target (update every C steps)

$$\boxed{\nabla_{\theta_t} L = E_{(s,a,r,s')\sim D}(r + \gamma \max_{a'} Q(s',a';\theta_{old}) - Q(s,a;\theta_t))\nabla_{\theta_t} Q(s,a;\theta_t)}$$

→ 不更新. C 步- 替换

fixed, old Q-target

**The Deep Q-Learning Algorithm(P27)**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$ 记忆体
Initialize action-value function $Q$ with random weights $\theta$
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**

→加一定的bias 来探索 避免局部最优解

        With probability $\epsilon$ select a random action $a_t$     **Sampling**
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$     **Training**
    **end for**
**end for**

# PART2

## Policy Gradient: Key Idea(P31)

**Policy Gradient**: directly optimize the policy $\pi(s)$



$P(a_1|s) = 0.9$
$P(a_2|s) = 0.1$
$P(a_3|s) = 0$

$\Rightarrow \pi^*(s) \sim P(a|s) = a_1$

State, s      **DNN**