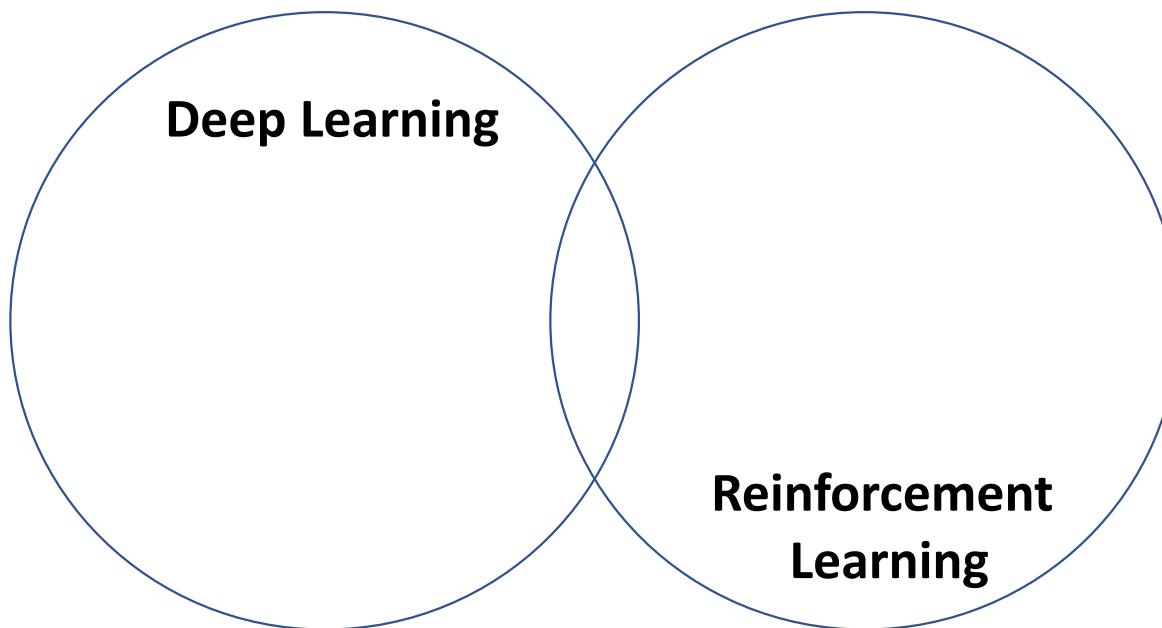


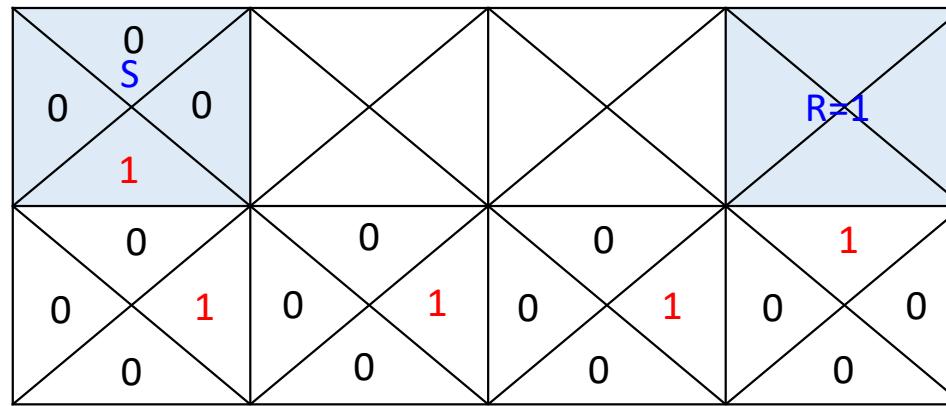
# Reinforcement Learning

# **Reinforcement Learning**



# Q-Learning

Q-Learning: Greedy action in random direction of Q-map



# Q-Learning: Exploration

- Exploitation

- ✓ Greedy action using Q-map

- $\epsilon$ -Greedy

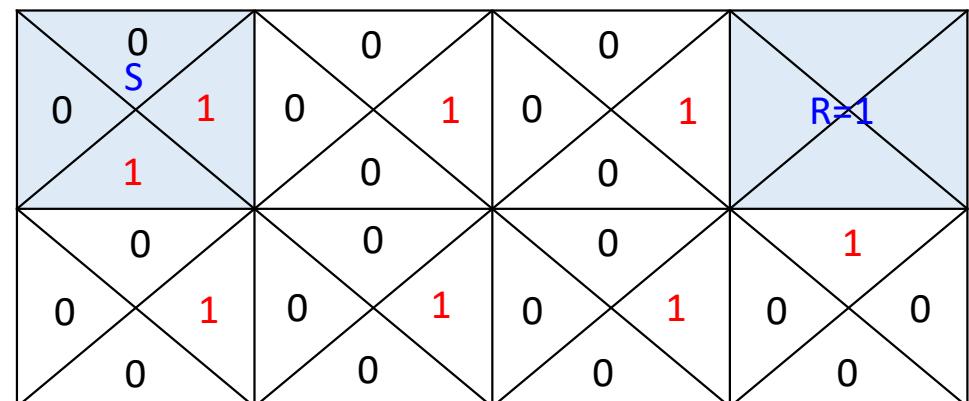
- ✓  $\epsilon(0,1)$ : random selection in direction + Greedy action

- Exploration

- ✓ Find new paths and goals

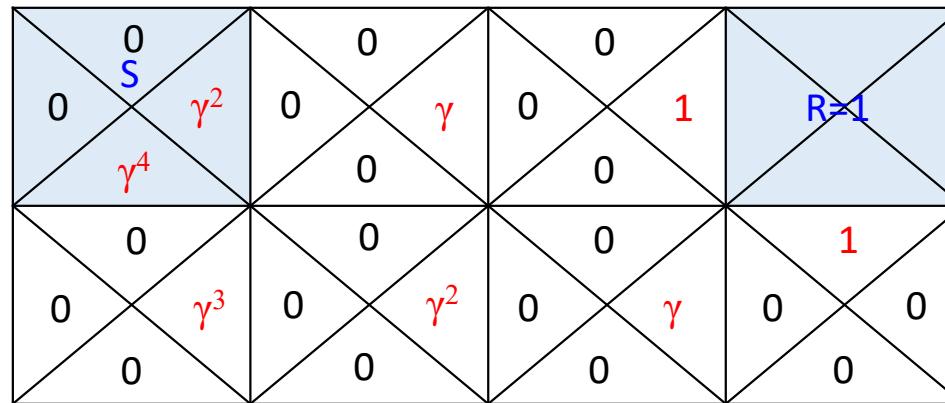
- (Decaying) $\epsilon$ -Greedy

- ✓ Variable  $\epsilon$  (e.g.,  $0.9 \rightarrow 0$ )



# Q-Learning: Discount factor

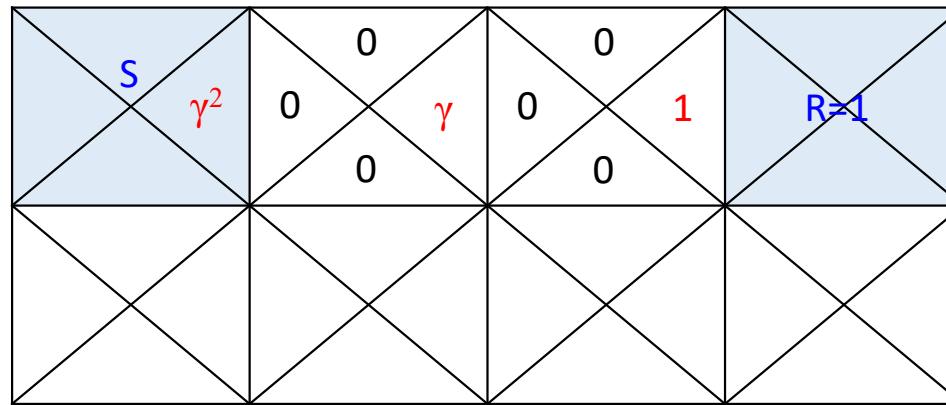
- Discount factor:  $\gamma$  (0, 1)



- ✓ Find efficient paths
- ✓ Meaning: reward for future state

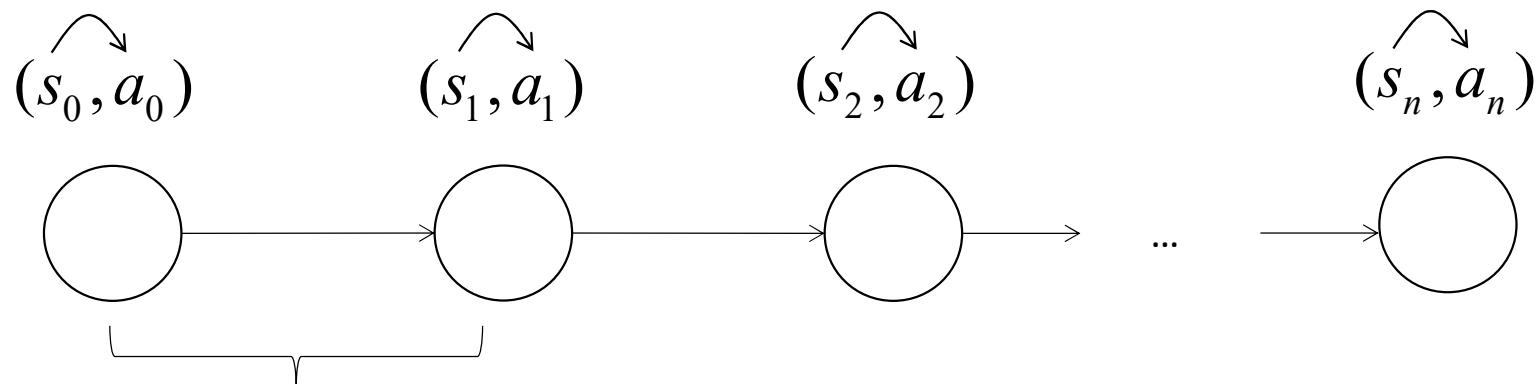
# Q-Learning: Q-update

- $Q(s_t, a_t)$



$$Q(s_t, a_t) \leftarrow \underbrace{(1 - \alpha)Q(s_t, a_t)}_{\text{Existing value}} + \underbrace{\alpha \left( R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \right)}_{\text{Newly updated value}}$$

# Markov Decision Process (MDP)



$P(a_1 | s_0, a_0, s_1)$  Prob. of action at state  $s_1$  (Policy)

$s_1$  has the information of  $s_0$  and  $a_0 = P(a_1 | s_1)$

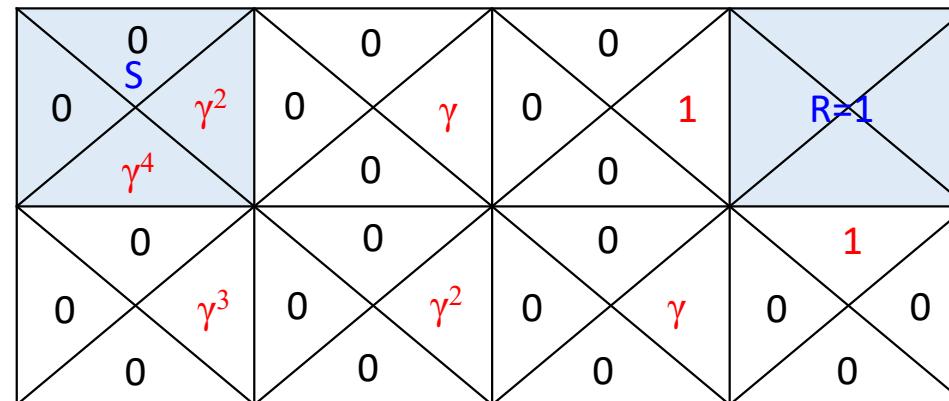
$P(s_2 | s_0, a_0, s_1, a_1) = P(s_2 | s_1, a_1)$

Prob. of state transition from  $s_1$  to  $s_2$

# Markov Decision Process (MDP)

- Goal of Reinforcement Learning
  - ✓ Maximize expected reward or return
    - Reward: result in action  $a_t$

$$G_t \cong R_t + \gamma R_{t+1} + r^2 R_{t+2} + \dots$$



# Markov Decision Process (MDP)

- State value function:
  - ✓ Expected return value in current state
- Action value function:
  - ✓ Expected return from current action
- Optimal policy
  - ✓ Maximize state value function

# Markov Decision Process (MDP)

- State value function  $V(s_t)$

$$G_t \cong R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

Expected Return

$$\Rightarrow V(s_t) \cong \underbrace{\int_{a_t:a_\infty} G_t P(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t : a_\infty}_{\text{For whole actions and states}}$$

Expected value for whole actions and states in  $S_t$

$$*E[f(x)] = \int f(x)p(x)dx$$

# Markov Decision Process (MDP)

- Action value function  $Q(s_t, a_t)$

$$G_t \cong R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

Expected Return

$$\rightarrow Q(s_t, a_t) \cong \int_{s_{t+1}:a_\infty} G_t P(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} : a_\infty$$

For whole actions and states

Expected value for whole actions and states from  $a_t$  in  $s_t$

# Markov Decision Process (MDP)

- Optimal policy: Maximize  $V(s_t)$

$$V(s_t) \cong \int_{a_t : a_\infty} G_t P(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t : a_\infty$$

Optimal policy

$$\left[ \begin{array}{l} P(a_t | s_t) \\ P(a_{t+1} | s_{t+1}) \\ \vdots \\ P(a_\infty | s_\infty) \end{array} \right]$$

# Bellman Equation

- Use Bayesian rule  $P(x, y) = P(x | y)P(y)$

$$P(x, y | z) = P(x | y, z)P(y | z)$$

$$V(s_t) \cong \int_{a_t : a_\infty} G_t P(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t : a_\infty$$

$\xrightarrow{\hspace{1cm}}$   $P(s_{t+1}, a_{t+1}, \dots | s_t, a_t)P(a_t | s_t)$

$$\begin{aligned} V(s_t) &\cong \int_{a_t : a_\infty} G_t P(s_{t+1}, a_{t+1}, \dots | s_t, a_t) P(a_t | s_t) da_t : a_\infty \\ &= \int_{a_t} \boxed{\int_{s_{t+1} : a_\infty} G_t P(s_{t+1}, a_{t+1}, \dots | s_t, a_t) ds_{t+1} : a_\infty} P(a_t | s_t) da_t \\ &= \int_{a_t} Q(s_t, a_t) P(a_t | s_t) da_t \quad \text{State value is the mean Q-values at the given state} \end{aligned}$$

# Bellman Equation

$$V(s_t) \cong \int_{a_t : a_\infty} G_t \underbrace{P(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}_{\rightarrow P(s_{t+1}, a_{t+1}, \dots | s_t, a_t) P(a_t | s_t)} da_t : a_\infty$$
$$\rightarrow P(a_{t+1}, \dots | s_t, a_t, s_{t+1}) P(a_t, s_{t+1} | s_t)$$
$$= P(a_{t+1}, \dots | s_{t+1}) P(a_t, s_{t+1} | s_t)$$

# Bellman Equation

$$G_t \cong R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = R_t + \gamma G_{t+1}$$

$$\begin{aligned} V(s_t) &\cong \int_{a_t, s_{t+1}} \left[ \int_{a_{t+1}:a_\infty} (R_t + \gamma G_{t+1}) P(a_{t+1}, \dots | s_{t+1}) da_{t+1} : a_\infty \right] P(a_t, s_{t+1} | s_t) da_t, s_{t+1} \\ &= \int_{a_t, s_{t+1}} (R_t + \gamma V(s_{t+1})) \underbrace{P(a_t, s_{t+1} | s_t)}_{\text{Policy}} da_t, s_{t+1} \\ &= \int_{a_t, s_{t+1}} (R_t + \gamma V(s_{t+1})) \underbrace{P(s_{t+1} | s_t, a_t)}_{\text{Transition}} \underbrace{P(a_t | s_t)}_{\text{Policy}} da_t, s_{t+1} \xrightarrow{\text{Bellman eq.}} \end{aligned}$$

# Bellman Equation

$$Q(s_t, a_t) \cong \int_{s_{t+1}:a_\infty} G_t P(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} : a_\infty$$

$\xrightarrow{\hspace{1cm}}$   $P(a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t, s_{t+1}) P(s_{t+1} | s_t, a_t)$

$$= P(a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_{t+1}) P(s_{t+1} | s_t, a_t)$$

$$G_t \cong R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = R_t + \gamma G_{t+1}$$

$$Q(s_t, a_t) \cong \int_{s_{t+1}} \left[ \int_{a_{t+1}:a_\infty} (R_t + \gamma G_{t+1}) P(a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_{t+1}) da_{t+1} : a_\infty \right] P(s_{t+1} | s_t, a_t) ds_{t+1}$$
$$= \int_{s_{t+1}} (R_t + \gamma V(s_{t+1})) P(s_{t+1} | s_t, a_t) ds_{t+1}$$

# Bellman Equation

$$Q(s_t, a_t) \equiv \int_{s_{t+1}:a_\infty} G_t P(a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t, s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1} : a_\infty$$

$\xrightarrow{\hspace{1cm}}$   $P(s_{t+2}, a_{t+2}, \dots | s_{t+1}, a_{t+1}) P(s_{t+1}, a_{t+1} | s_t, a_t)$

$$\begin{aligned} Q(s_t, a_t) &\equiv \int_{s_{t+1}, a_{t+1}} \left[ \int_{s_{t+2}:a_\infty} (R_t + \gamma G_{t+1}) P(s_{t+2}, a_{t+2}, \dots | s_{t+1}, a_{t+1}) ds_{t+2} : a_\infty \right] P(s_{t+1}, a_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \\ &= \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) \underbrace{P(s_{t+1}, a_{t+1} | s_t, a_t)}_{\text{Policy}} ds_{t+1}, a_{t+1} \\ &= \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_t, a_t, s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \\ &= \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) \underbrace{P(a_{t+1} | s_{t+1})}_{\text{Policy}} \underbrace{P(s_{t+1} | s_t, a_t)}_{\text{Transition}} ds_{t+1}, a_{t+1} \xrightarrow{\hspace{1cm}} \text{Bellman eq.} \end{aligned}$$

# Optimal Policy

- Maximize the expected return at the current state
- Maximize the state value function

$$\begin{aligned} V(s_t) &\equiv \int_{a_t : a_\infty} G_t P(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t : a_\infty \\ &= \int_{a_t} Q(s_t, a_t) P(a_t | s_t) da_t \end{aligned}$$

$$\arg \max_{P(a_t | s_t)} \int_{a_t} Q(s_t, a_t) P(a_t | s_t) da_t \iff \int_{a_t} \underline{Q^*(s_t, a_t)} P(a_t | s_t) da_t$$

$$a_t^* \equiv \arg \max_{a_t} Q^*(s_t, a_t) \Rightarrow \text{Greedy policy} \iff \varepsilon\text{-Greedy for random direction}$$

$$P^*(a_t | s_t) = \delta(a_t - a_t^*) \text{ Find } a_t \text{ to maximize } Q^* \Rightarrow \text{Choose } a_t$$

# Monte-Carlo Method

- How to obtain  $Q^*$ ?
  - ✓ Training means making  $Q^*$

- Expectation value  $E[x] = \int_x xP(x)dx \approx \frac{1}{N} \sum_{i=1}^N x_i$

$$Q(s_t, a_t) \cong \int_{s_{t+1}:a_\infty} G_t P(s_{t+1} : a_\infty | s_t, a_t) ds_{t+1} : a_\infty = \frac{1}{N} \sum_{i=1}^N G_t^{(i)}$$

For many cases  $N$ , obtain  $Q$  and find  $a_t$

# Temporal Difference (TD)

- Action value function

$$\begin{aligned} Q(s_t, a_t) &\cong \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(s_{t+1} | s_t, a_t) P(a_{t+1} | s_{t+1}) ds_{t+1}, a_{t+1} \\ &\approx \frac{1}{N} \sum_{i=1}^N (R_t^{(i)} + \gamma Q(s_{t+1}^{(i)}, a_{t+1}^{(i)})) \cong \bar{Q}_N \end{aligned}$$

- 1-step TD (Incremental Monte-Carlo Updates)

$$\begin{aligned} \bar{Q}_N &= \frac{1}{N} \left( \bar{Q}_{N-1}(N-1) + R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right) \\ &= \underline{\bar{Q}_{N-1}} + \frac{1}{N} \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) - \underline{\bar{Q}_{N-1}} \right) \end{aligned}$$

# Incremental Monte-Carlo Updates

- Incremental Mete-Carlo Updates

$$\begin{aligned}\bar{Q}_N &= \frac{1}{N} \left( \bar{Q}_{N-1}(N-1) + R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right) \\ &= \bar{Q}_{N-1} + \boxed{\frac{1}{N}} \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) - \bar{Q}_{N-1} \right) \\ &= (1-\alpha)\bar{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right) \\ &\quad \text{Current sample} \\ &\quad (\text{i.e., TD target})\end{aligned}$$

- SARSA algorithm

$$\bar{Q}_N = (1-\alpha)\bar{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right)$$

Reward

Current State & Action      Next State      Next Action

# Temporal Difference (TD)

- SARSA
  - ✓ On-Policy
    - Behavior policy == Target policy
  
- Q-Learning
  - ✓ Off-Policy
    - Behavior policy != Target policy

# Behavior policy & Target policy

- Behavior policy & Target policy

$$Q(s_t, a_t) \cong \underbrace{\int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}}$$

Target policy

Given by Behavior policy  
 $P(a_t | s_t)$

$$\approx \overline{Q}_N = \underbrace{(1 - \alpha) \overline{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right)}_{\text{TD target} \\\text{(i.e., sample)}}$$

# Off-Policy

- Divide Target policy and Behavior policy

✓ Why?

- Can use other target policy by human or other agents
- Can obtain a sample (i.e., TD target) by the optimal policy (i.e., greedy action) during exploring
- Re-judgement (i.e., re-use) for actions is possible

$$Q(s_t, a_t) \cong \underbrace{\int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}}_{\text{Target policy}}$$

Behavior policy:  
 $p(a_t | s_t)$

$$\approx \bar{Q}_N = (1 - \alpha) \bar{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right)$$

# Q-Learning (revisit)

- Target policy: greedy
- Behavior policy:  $\epsilon$ -greedy

$$Q(s_t, a_t) \cong \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}$$

$$\approx \bar{Q}_N = (1 - \alpha) \bar{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right)$$

Target :  $P(a_{t+1} | s_{t+1}) = \delta(a_{t+1} - a_{t+1}^*)$ ,  $a_{t+1}^* = \arg \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$

Behavior :  $\int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) \delta(a_{t+1} - a_{t+1}^*) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}$

$$= \int_{s_{t+1}} (R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})) P(s_{t+1} | s_t, a_t) ds_{t+1}$$

# Q-Learning (revisit)

$$\begin{aligned} & \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) \delta(a_{t+1} - a_{t+1}^*) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \\ &= \int_{s_{t+1}} (R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})) P(s_{t+1} | s_t, a_t) ds_{t+1} \\ &\Rightarrow R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) : \text{TD-target} \end{aligned}$$

## 2Step-TD

- 1-step TD

$$\begin{aligned} Q(s_t, a_t) &\cong \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \\ &\approx \bar{Q}_N = (1 - \alpha) \bar{Q}_{N-1} + \alpha \left( R_t^{(N)} + \gamma Q(s_{t+1}^{(N)}, a_{t+1}^{(N)}) \right) \end{aligned}$$

- 2-step TD

$$\begin{aligned} Q(s_t, a_t) &\cong \\ &\int_{s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}} (R_t + \gamma R_t + \gamma^2 Q(s_{t+2}, a_{t+2})) P(a_{t+2} | s_{t+2}) P(s_{t+2} | s_{t+1}, a_{t+1}) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}, s_{t+2}, a_{t+2} \end{aligned}$$

## 2Step-TD

- Off-policy with importance sampling

$$\int_x xp(x)dx \approx \frac{1}{N} \sum_{i=1}^N x_i \quad x_i \sim p(x)$$
$$= \int_x x \frac{p(x)}{q(x)} q(x)dx \approx \frac{1}{N} \sum_{i=1}^N x_i \boxed{\frac{p(x_i)}{q(x_i)}} \quad x_i \sim q(x)$$

Compensation value  $\leftarrow$  importance rate

## 2Step-TD

- TD-target in 2-step TD

$$Q(s_t, a_t) \cong \int_{s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}} (R_t + \gamma R_t + \gamma^2 Q(s_{t+2}, a_{t+2})) P(a_{t+2} | s_{t+2}) P(s_{t+2} | s_{t+1}, a_{t+1}) \frac{P(a_{t+1} | s_{t+1})}{q(a_{t+1} | s_{t+1})} P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}$$

Not Use : We can obtain similar results

$$\text{TD-target : } \frac{p(a_{t+1}^{(N)} | s_{t+1})}{q(a_{t+1}^{(N)} | s_{t+1})} \left( R_t^{(N)} + \gamma R_{t+1}^{(N)} + \gamma^2 \max_{a_{t+2}} Q(s_{t+2}^{(N)}, a_{t+2}) \right)$$

$a_t^{(N)} \sim q$     $a_{t+1}^{(N)} \sim q$     $a_{t+2}^{(N)} \sim p$

# Q-learning Code

```
import numpy as np

n_states = 12 # Number of states in the grid world
n_actions = 4 # Number of possible actions (up, down, left, right)
goal_state = 11 # Goal state

Q_table = np.zeros((n_states, n_actions))

learning_rate = 0.8
discount_factor = 0.95
exploration_prob = 0.2
epochs = 100
```

# Q-learning Code

```
for epoch in range(epochs):
    current_state = np.random.randint(0, n_states) # Start from a random state

    while current_state != goal_state:
        # Choose action with epsilon-greedy strategy
        if np.random.rand() < exploration_prob:
            action = np.random.randint(0, n_actions) # Explore
        else:
            action = np.argmax(Q_table[current_state]) # Exploit

        # Simulate the environment (move to the next state)
        # For simplicity, move to the next state
        next_state = (current_state + 1) % n_states

        # Define a simple reward function (1 if the goal state is reached, 0 otherwise)
        reward = 1 if next_state == goal_state else 0

        # Update Q-value using the Q-learning update rule
        Q_table[current_state, action] += learning_rate * (reward +
                                                          discount_factor * np.max(Q_table[next_state]) - Q_table[current_state, action])

        current_state = next_state # Move to the next state

    # After training, the Q-table represents the learned Q-values
    print("Learned Q-table:")
    print(Q_table)
```

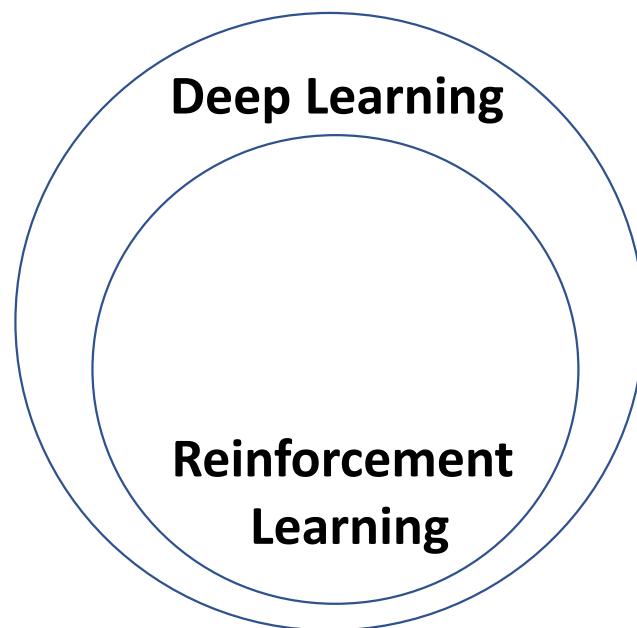
# Q-learning Results

Learned Q-table:

```
[[0.37616543 0.47823792 0.59873694 0.          ]
 [0.63024941 0.          0.60503942 0.          ]
 [0.66342043 0.53073304 0.65809562 0.50961706]
 [0.6983373  0.55866984 0.6927506  0.6927506 ]
 [0.73509189 0.72921116 0.73385176 0.7292102 ]
 [0.77378094 0.7428297  0.77251646 0.76759069]
 [0.81450625 0.80799003 0.81449582 0.781926  ]
 [0.857375   0.85710064 0.85595476 0.85736402]
 [0.9025     0.90244224 0.90244224 0.901056  ]
 [0.94999749 0.94998776 0.95       0.95       ]
 [1.          0.99999949 0.99999949 0.99999998]
 [0.          0.          0.          0.          ]]
```

# Deep Q Network (DQN)

**Q-Learning +DNN**



# Q-Learning

$$Q(s_t, a_t) \cong \int_{s_{t+1}, a_{t+1}} (R_t + \gamma Q(s_{t+1}, a_{t+1})) P(a_{t+1} | s_{t+1}) P(s_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1}$$

Target policy    Transition

$$P(a_{t+1} | s_{t+1}) = \delta(a_{t+1} - a_{t+1}^*), \quad a_{t+1}^* = \arg \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

$$= \int_{s_{t+1}} (R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})) P(s_{t+1} | s_t, a_t) ds_{t+1}$$

TD target

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left( R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \right)$$

# Q-Learning

- Bellman equation as an iterative update

$$Q(s_t, a_t) \leftarrow \mathbb{E}_{s_{t+1} \sim P(s_{t+1}|s_t, a_t)} \left( R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \right)$$

- Q will converge to  $Q^*$

$$\text{target}(s_{t+1}) = R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (\text{target}(s_{t+1}) - Q(s_t, a_t))$$

✓ Q-learning converges to optimal policy if all state-action pairs seen frequently enough

- $\epsilon$ -Greedy as Behavior policy

# Problem in Q-learning

- Not scalable
  - ✓ Must compute Q-table for every state-action pair
- Use a function approximator to estimate  $Q(s,a)$ 
  - ✓ Learn about some small number of training states from experience
  - ✓ Generalize that experience to new, similar situations

# Deep Q Network (DQN)

- Regression with DNN
    - ✓ Find Q using Regression
      - Why?
      - To apply various states
    - ✓  $Q_w$  : Q with weight w; By regression

$$\vec{w} \leftarrow \vec{w} - \boxed{\alpha} \text{Grad}$$

Learning rate

Difference  
among samples

# Approximate Q-Learning

$$\text{target}(s_{t+1}) = R_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (\text{target}(s_{t+1}) - Q(s_t, a_t))$$

- Linear function approximator (i.e., Regression)

$$Q_w(s, a) = w_0 f_0(s, a) + w_1 f_1(s, a) + \dots + w_n f_n(s, a)$$

- Exact Q

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$$

- Approximate Q

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$$

# Linear function approximator

- Approximate Q  $w_i \leftarrow w_i + \alpha [difference] f_i(s, a)$

- Loss function

$$loss = \frac{1}{2}(y - \hat{y})^2$$

- Difference

$$y - \hat{y} = y - (w_0 f_0(s, a) + w_1 f_1(s, a) + \dots + w_n f_n(s, a)) = y - \sum_i w_i f_i(s, a)$$

$$\frac{\partial l(w)}{\partial w_k} = - \left( y - \sum_i w_i f_i(s, a) \right) f_k(s, a)$$

$$w_k \leftarrow w_i + \alpha \left( y - \sum_i w_i f_i(s, a) \right) f_k(s, a)$$

$$w_k \leftarrow w_i + \alpha [difference] f_k(s, a)$$

# Deep Q Network (DQN)

- Difference

$$difference = \left[ R(s_t) + \gamma \max_{a_{t+1}} Q_w(s_{t+1}, a_{t+1}) \right] - Q_w(s_t, a_t)$$

- Loss function

$$l(w) = \left[ \frac{1}{2} (\text{target}(s_{t+1}) - Q_w(s_t, a_t))^2 \right]$$

- TD update

$$\begin{aligned} w_{k+1} &\leftarrow w_k - \alpha \nabla_w l(w) \\ &\leftarrow w_k - \alpha \nabla_w \left[ \frac{1}{2} (\text{target}(s_{t+1}) - Q_w(s_t, a_t))^2 \right] \\ &\leftarrow w_k - \alpha [\text{target}(s_{t+1}) - Q_w(s_t, a_t)] \nabla_w Q_w(s_t, a_t) \end{aligned}$$

# Deep Q Network 2013

---

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

---

Initialize replay memory  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

**for** episode = 1,  $M$  **do**

    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$

**for**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$

        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$

        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3

**end for**

**end for**

---

# Deep Q Network 2015

**Algorithm 1:** deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$

**For** episode = 1,  $M$  **do**

    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$

        otherwise select  $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$

        Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$

        Every  $C$  steps reset  $\hat{Q} = Q$

**End For**

**End For**

# DQN Code

```
import gym
import collections
import random

import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

learning_rate = 0.0005
gamma = 0.9
buffer_limit = 10000
batch_size = 64

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

# DQN Code

```
class ReplayBuffer():
    def __init__(self):
        self.buffer = collections.deque(maxlen=buffer_limit)
    def put(self, transition):
        self.buffer.append(transition)
    def sample(self, n):
        mini_batch = random.sample(self.buffer, n)
        s_list, a_list, r_list, s_prime_list, done_mask_list = [], [], [], [], []

        for transition in mini_batch:
            s, a, r, s_prime, done_mask = transition
            s_list.append(s)
            a_list.append([a])
            r_list.append([r])
            s_prime_list.append(s_prime)
            done_mask_list.append([done_mask])

        return torch.tensor(s_list, dtype=torch.float), torch.tensor(a_list), torch.tensor(r_list),\
               torch.tensor(s_prime_list, dtype=torch.float), torch.tensor(done_mask_list)
    def size(self):
        return len(self.buffer)
```

# DQN Code

```
class Qnet(nn.Module):
    def __init__(self):
        super(Qnet, self).__init__()
        self.fc1 = nn.Linear(4, 256)
        self.fc2 = nn.Linear(256, 2)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = self.fc2(x)
        return x

    def sample_action(self, obs, epsilon):
        out = self.forward(obs)
        prob = random.random()
        if prob < epsilon:
            return random.randint(0,1)
        else:
            return out.argmax().item()
```



# DQN Code

```
def train(q, q_target, memory, optimizer):
    for i in range(10):
        s, a, r, s_prime, done_mask = memory.sample(batch_size)

        q_out = q(s)
        q_a = q_out.gather(1,a)
        max_q_prime = q(s_prime).max(1)[0].unsqueeze(1)
        target = r + gamma * max_q_prime * done_mask
        loss = F.smooth_l1_loss(q_a, target)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```



```

def main():
    env = gym.make('CartPole-v1')
    q = Qnet()
    q_target = Qnet()
    q_target.load_state_dict(q.state_dict())
    memory = ReplayBuffer()

    print_interval = 20
    score = 0.0
    optimizer = optim.Adam(q.parameters(), lr=learning_rate)

    for n_epi in range(1000):
        epsilon = max(0.01, 0.08 - 0.01*(n_epi/200)) # Linear annealing from 8% to 1%
        s, _ = env.reset()
        done = False

        while not done:
            a = q.sample_action(torch.from_numpy(s).float(), epsilon)
            s_prime, r, done, truncated, info = env.step(a)
            done_mask = 0.0 if done else 1.0
            memory.put((s, a, r, s_prime, done_mask))
            s = s_prime
            score += r
            if done:
                break

        if memory.size() > 2000:
            train(q, q_target, memory, optimizer)

        if n_epi%print_interval==0 and n_epi!=0:
            q_target.load_state_dict(q.state_dict())
            print("n_episode :{}, score : {:.1f}, n_buffer : {}, eps : {:.1f}%".format(
                n_epi, score/print_interval, memory.size(), epsilon*100))
            score = 0.0

```

# DQN - Results

```
n_episode :500, score : 58.4, n_buffer : 10000, eps : 5.5%
n_episode :520, score : 81.7, n_buffer : 10000, eps : 5.4%
n_episode :540, score : 93.5, n_buffer : 10000, eps : 5.3%
n_episode :560, score : 98.2, n_buffer : 10000, eps : 5.2%
n_episode :580, score : 94.5, n_buffer : 10000, eps : 5.1%
n_episode :600, score : 83.9, n_buffer : 10000, eps : 5.0%
n_episode :620, score : 86.5, n_buffer : 10000, eps : 4.9%
n_episode :640, score : 83.2, n_buffer : 10000, eps : 4.8%
n_episode :660, score : 85.5, n_buffer : 10000, eps : 4.7%
n_episode :680, score : 102.3, n_buffer : 10000, eps : 4.6%
n_episode :700, score : 125.3, n_buffer : 10000, eps : 4.5%
n_episode :720, score : 128.8, n_buffer : 10000, eps : 4.4%
n_episode :740, score : 179.3, n_buffer : 10000, eps : 4.3%
n_episode :760, score : 162.4, n_buffer : 10000, eps : 4.2%
n_episode :780, score : 168.2, n_buffer : 10000, eps : 4.1%
n_episode :800, score : 135.7, n_buffer : 10000, eps : 4.0%
n_episode :820, score : 133.3, n_buffer : 10000, eps : 3.9%
n_episode :840, score : 134.6, n_buffer : 10000, eps : 3.8%
n_episode :860, score : 145.1, n_buffer : 10000, eps : 3.7%
n_episode :880, score : 170.6, n_buffer : 10000, eps : 3.6%
n_episode :900, score : 199.9, n_buffer : 10000, eps : 3.5%
n_episode :920, score : 193.4, n_buffer : 10000, eps : 3.4%
n_episode :940, score : 186.1, n_buffer : 10000, eps : 3.3%
n_episode :960, score : 169.6, n_buffer : 10000, eps : 3.2%
n_episode :980, score : 214.7, n_buffer : 10000, eps : 3.1%
Complete
```

E N D