

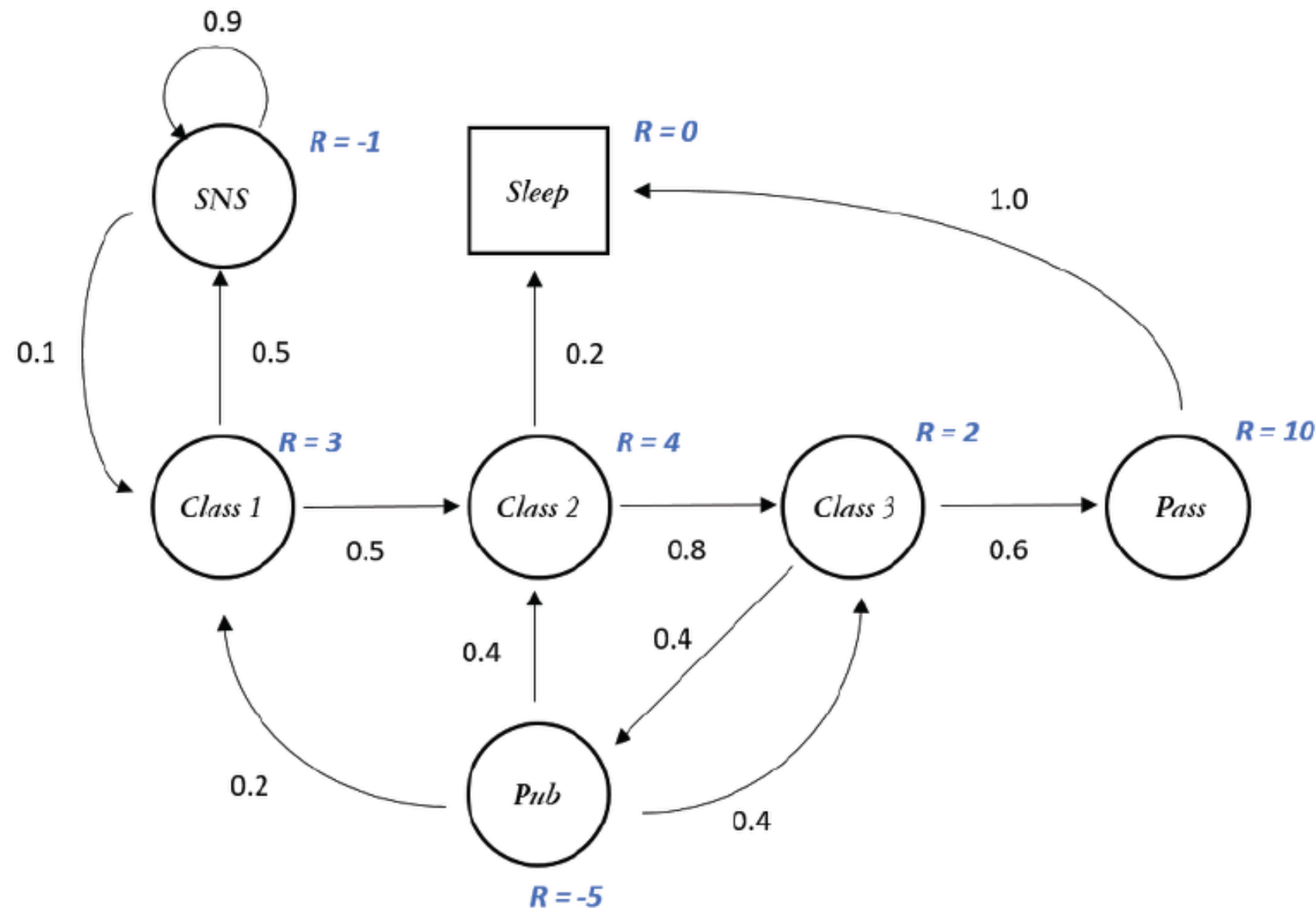
Introduction to RL – Part 2

TA. Bogyeeong Suh

Recap of Part 1

Markov Decision Process

- The process of RL is usually given as **Markov Decision Processes (MDPs)**, which is a mathematical framework used for modeling decision making in situations where the outcomes are partly random and partly under the control of a decision maker.



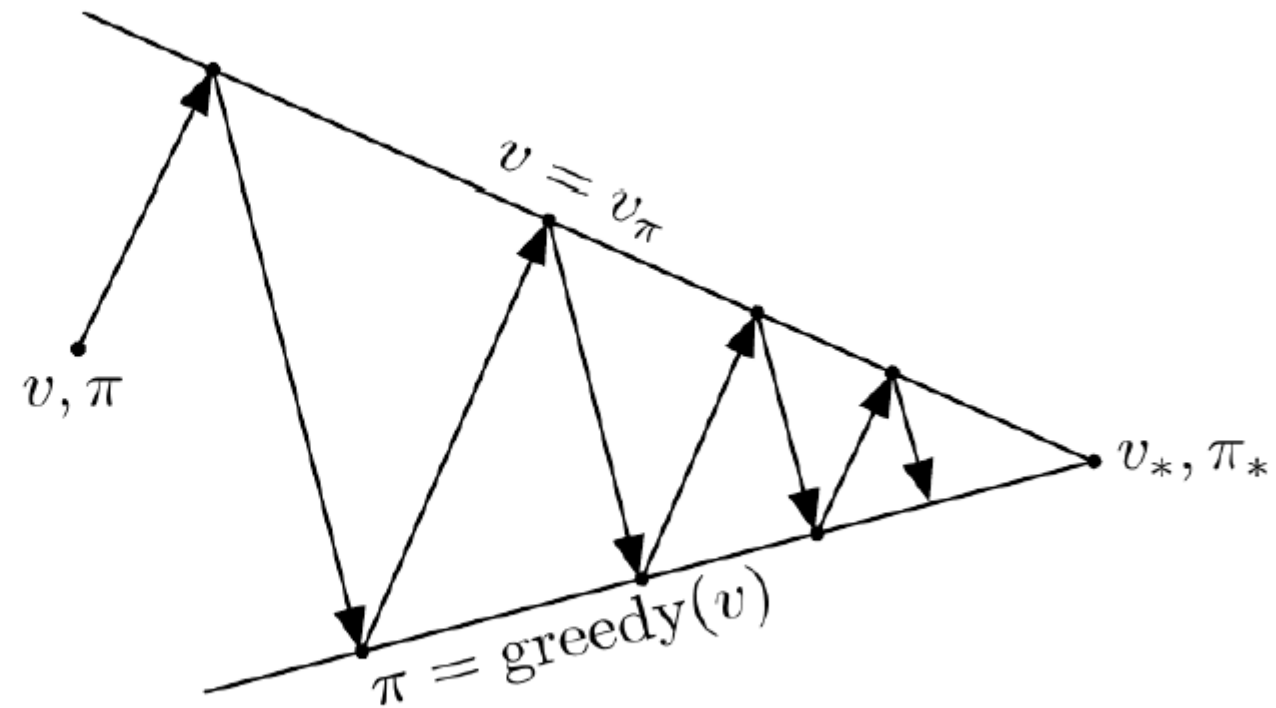
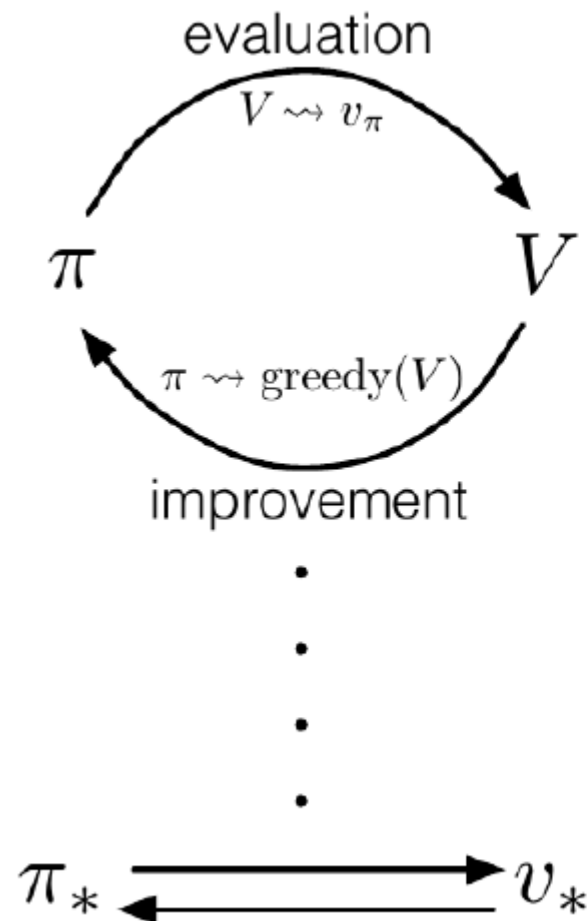
Markov Reward Process (MRP)

Recap of Part 1

Policy Iteration

1) Policy Evaluation

-Given an arbitrary policy, we calculate value function V for **all states** under this policy. We iterate through each state and update V until it converges. → We find true value function with iteration



Recap of Part 1

Value Iteration

$k = 0$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

at (1,2) state : state 1

Up $V_1(s) = -1 + 0$
 Down $V_1(s) = -1 + 0$
 Left $V_1(s) = -1 + 0$
 Right $V_1(s) = -1 + 0$

$$\therefore V_1(1) = \max V_1(s) = -1$$

$k = 1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

at (1,2) state : state 1

Up $V_2(s) = -1 + (-1)$
 Down $V_2(s) = -1 + (-1)$
 Left $V_2(s) = -1 + (-1)$
 Right $V_2(s) = -1 + (-1)$

$$\therefore V_2(s) = \max V_2(s) = -1$$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

at (1,3) state : state 2

Up $V_2(s) = -1 + (-1)$
 Down $V_2(s) = -1 + (-1)$
 Left $V_2(s) = -1 + (-1)$
 Right $V_2(s) = -1 + (-1)$

$$\therefore V_2(s) = \max V_2(s) = -2$$

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-1
-2	-2	-1	0

$k = 2$

0	-1	-2	-3
-1	-2	-3	-2
-2	-3	-2	-1
-3	-2	-1	0

$k = 3$

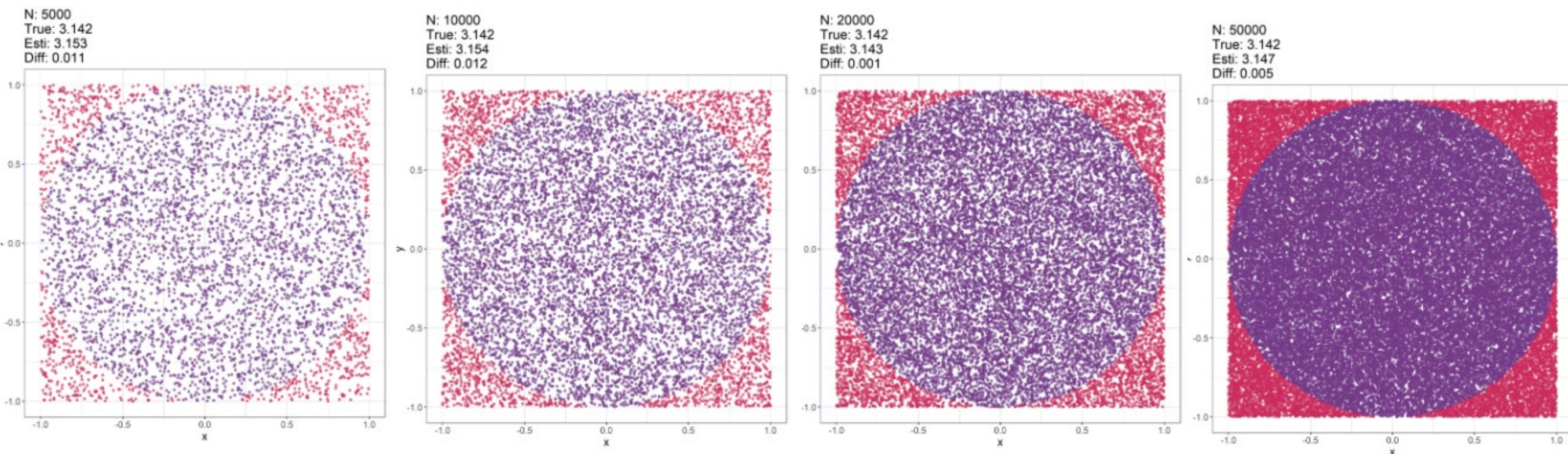
0	-1	-2	-3
-1	-2	-3	-2
-2	-3	-2	-1
-3	-2	-1	0

$k = \infty$

Monte Carlo method

Monte Carlo method

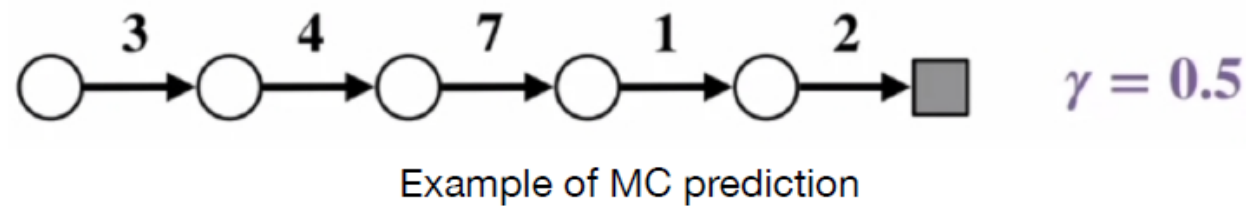
- Generally, MC is a method that relies on **repeated random sampling** to obtain numerical results
- MC methods learn directly from episodes of **experience**
- In MC, we learn from *complete* episodes: no bootstrapping
- MC is *model-free*: no knowledge of MDP transitions / rewards



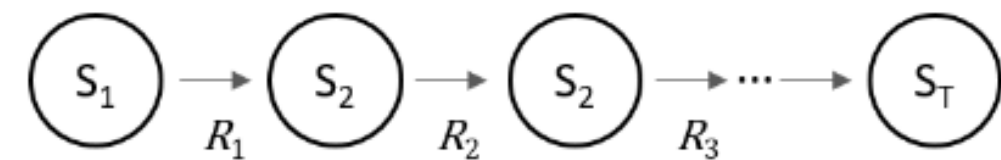
e.g. Finding the area of a circle with radius=1 using random samples of (x,y)

Monte Carlo method

Monte Carlo method



$$\begin{array}{ll}
 G_0 = R_1 + \gamma G_1 & G_5 = 0 \\
 G_1 = R_2 + \gamma G_2 & G_4 = R_5 + \gamma G_5 = 2 + 0.5 * 0 = 2 \\
 G_2 = R_3 + \gamma G_3 & G_3 = R_4 + \gamma G_4 = 1 + 0.5 * 2 = 2 \\
 G_3 = R_4 + \gamma G_4 & \rightarrow G_2 = R_3 + \gamma G_3 = 7 + 0.5 * 2 = 8 \\
 G_4 = R_5 + \gamma G_5 & G_1 = R_2 + \gamma G_2 = 4 + 0.5 * 8 = 8 \\
 G_5 = 0 & G_0 = R_1 + \gamma G_1 = 3 + 0.5 * 8 = 7
 \end{array}$$



$$\begin{aligned}
 G(s_1) &= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \\
 G(s_2) &= R_2 + \gamma^2 R_3 + \dots \\
 G(s_3) &= R_3 + \dots
 \end{aligned}$$

<Fig 2. The returns of each states in 1 episode>

- Goal: learn a policy that maximizes the total cumulative reward it receives over time
- Transition probability: $p(s', r | s, \pi(s))$
- State Value function: $V_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$ → **Average values of returns from episodes**
- Action value function: $q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$

Monte Carlo method

Monte Carlo method

1. Initialization

Initialize:

$\pi(s) \in A(s), \forall s \in \mathbb{S}$

$Q(s, a) \in \mathbb{R}, \forall s \in \mathbb{S}, \forall a \in A(s)$

$Returns(s, a) \leftarrow \text{empty list}, \forall s \in \mathbb{S}, \forall a \in A(s)$

2. Monte Carlo

Loop forever(for each episode):

Choose $S_0 \in \mathbb{S}, A_0 \in A(S_0)$, randomly such that all pairs have probability > 0
 Generate an episode from S_0, A_0 following $\pi: S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T - 1, T - 2, \dots, 0$

$G \leftarrow \gamma G + R_{t+1}$

Unless the S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{Average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

- Using randomly generated episodes, calculate action value function from cumulative reward G

$$q_{\pi}(s, a) = \frac{1}{N(s)} \sum_{i=1}^{N(s)} G_i(s)$$

- $N(s)$: number of times we visited state s we from start to end in total episodes
- $G_i(s)$: return of state s in episode i

Comparison of DP and MC

- **Dynamic Programming**

- **Pros: bootstrapping** (update value estimates based on other learned estimates, without waiting for a final outcome)
- **Cons: Model-Based** (requires complete and accurate model of the environment)

- **Monte Carlo method**

- **Pros: Model-free** (learn directly from experience)
- **Cons:** applicable to **episodic tasks** (require episodes to end to calculate returns, making them unsuitable for continuous tasks)

➔ **Temporal Difference Learning (TD) takes advantages from DP and MC; It is model-free** (uses random sampling), which learns directly from raw experience without a model of the environment, and **bootstraps** by updating estimates based in part on other learned estimates, without waiting for a final outcome.

Temporal Difference Learning

Temporal Difference Learning

Value function

$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$

Using Monte Carlo method

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

The episode should be finished to get the return G_t

Total expected return G_t

$$\begin{aligned} G_t &= \sum_{i=t+1}^T \gamma^{i-t-1} R_i \\ &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

State value function

$$\begin{aligned} V_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= R_{t+1} + \gamma v_{\pi}(S_{t+1}) \end{aligned}$$

Temporal Difference (TD) Learning

$$V(S_t) \leftarrow V(S_t) + \alpha[\underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{TD Target}} - \underbrace{V(S_t)}_{\text{TD Error}}]$$

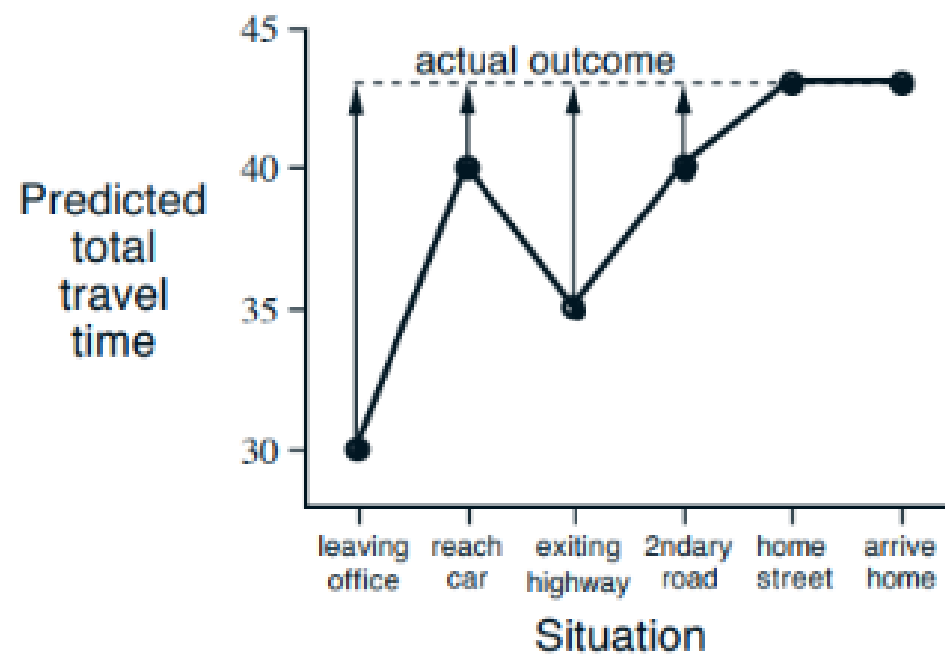
Estimated value of the current state

- It updates the value of the current state towards the estimated return
- TD error: measures the difference between the predicted value of the current state and the observed reward plus the estimated value of the next state

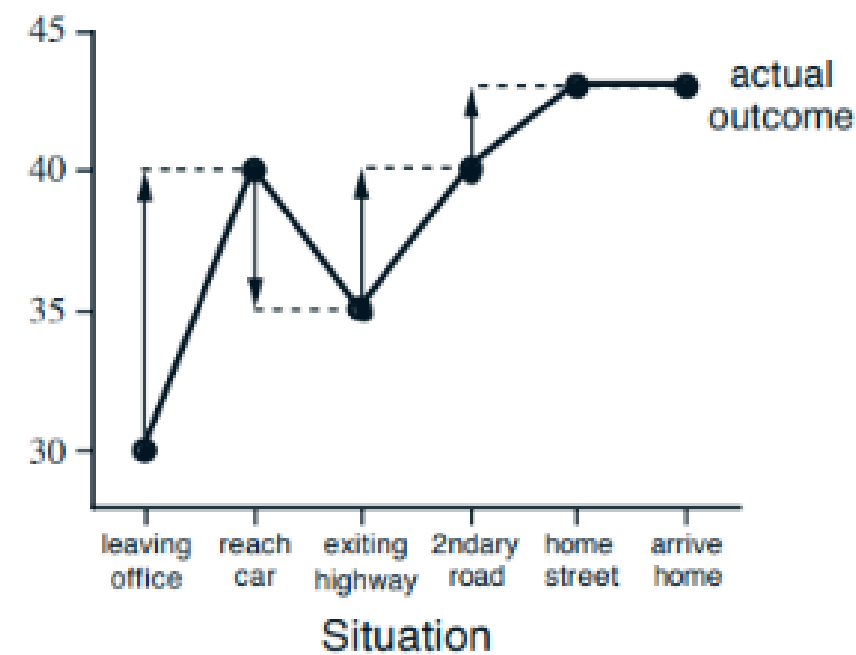
Temporal Difference Learning

Temporal Difference Learning

<i>State</i>	<i>Elapsed Time (minutes)</i>	<i>Predicted Time to Go</i>	<i>Predicted Total Time</i>
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43



Monte-Carlo method



Temporal Difference learning

SARSA

SARSA

-On-policy TD control

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \underbrace{(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))}_{\text{TD Error}}$$

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

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

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

Loop for each episode:

Initialize S

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

Loop for each step of episode:

Take action A , observe R, S'

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

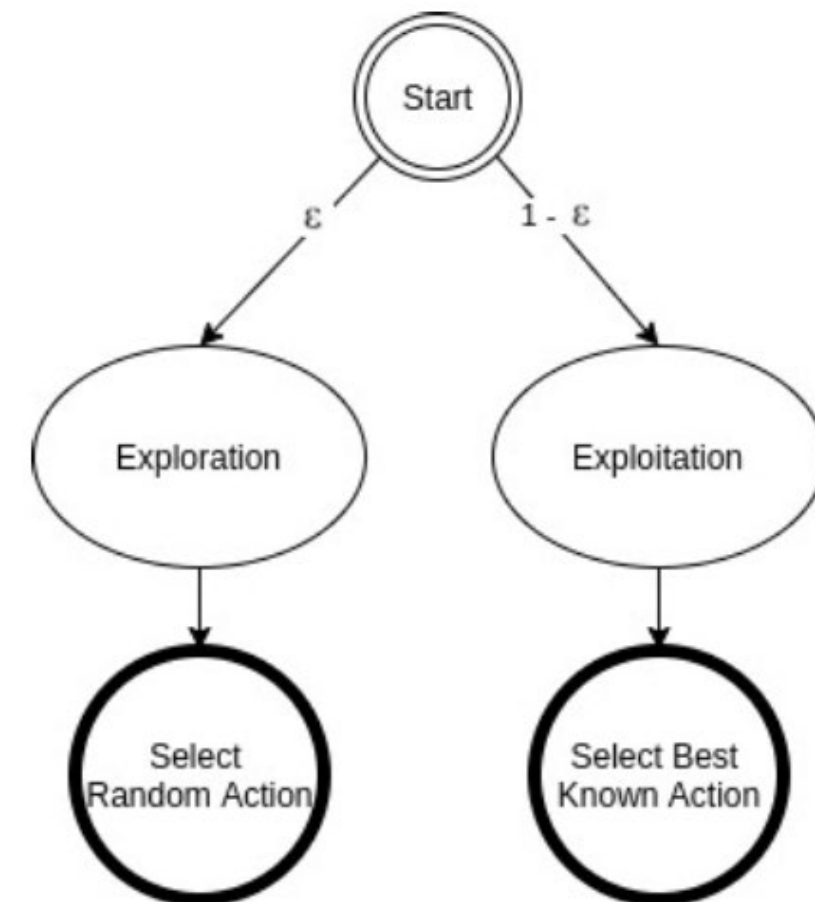
$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until S is terminal

ϵ -greedy

$$A \leftarrow \begin{cases} \arg \max_a Q(S, A) & \text{with probability } 1 - \epsilon \\ \text{a random action} & \text{with probability } \epsilon \end{cases}$$



ϵ -greedy

Q-learning

Q-learning

-Off-policy TD control

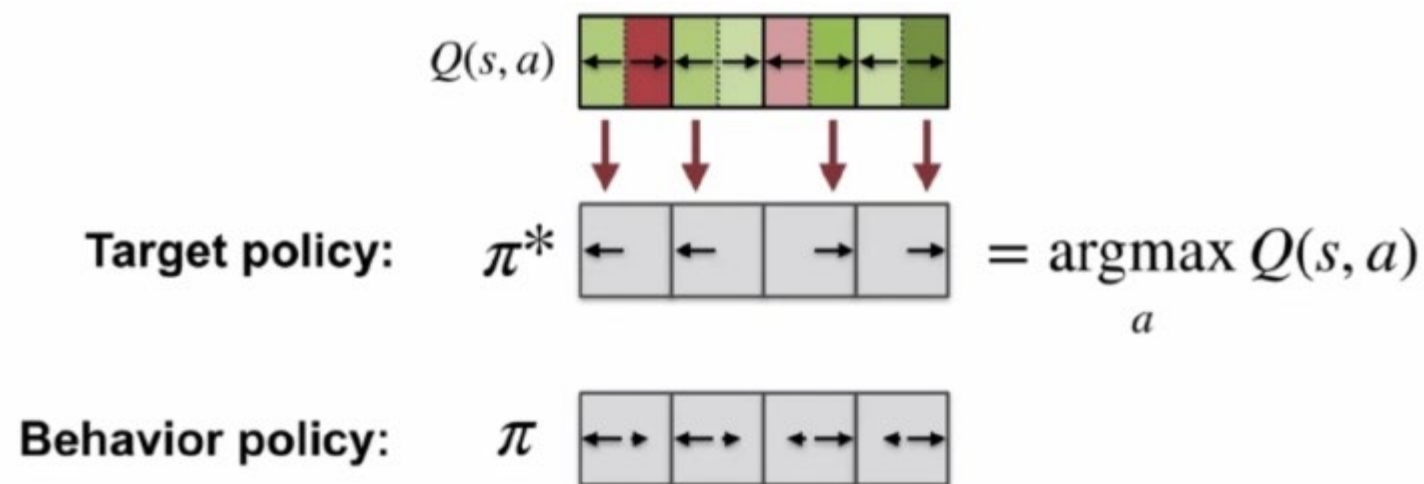
Update rule in SARSA

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

Update rule in Q-learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \underbrace{\max_{a'} Q(S_{t+1}, a')}_{\text{From } \pi_*} - Q(S_t, A_t))$$

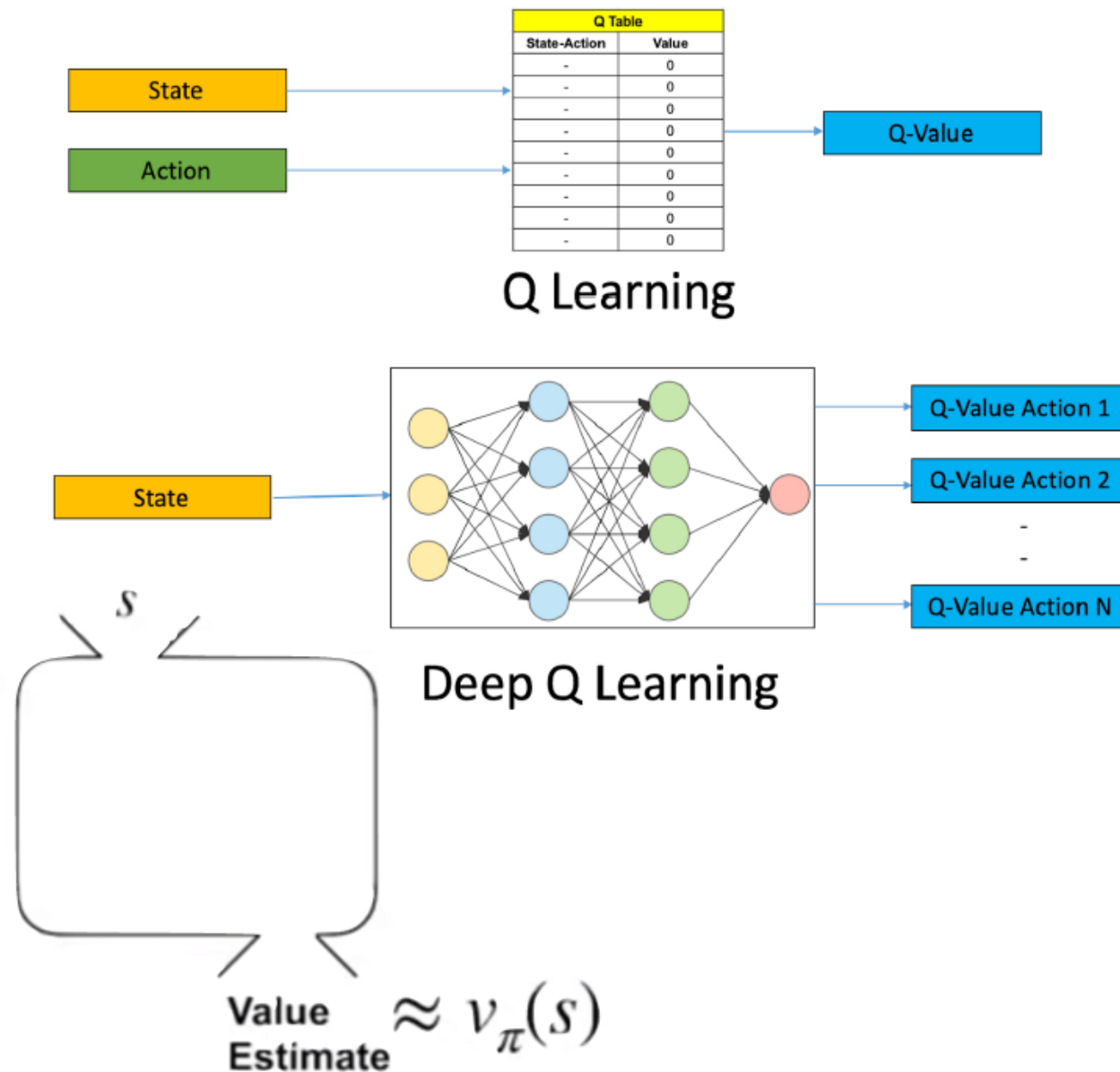
*Target policy π_**



Target and behavior policies

Deep Q-learning

Deep Q-learning



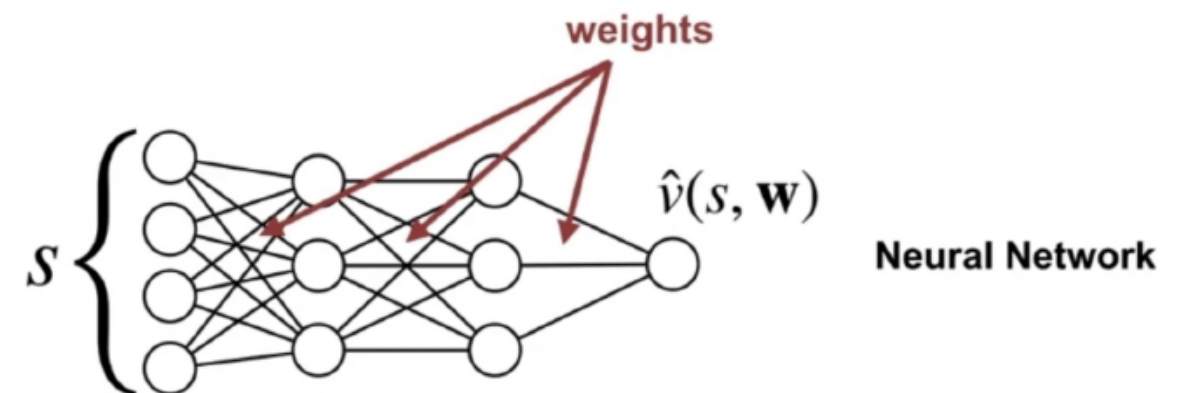
Parameterized value function

$$\hat{v}(s, \mathbf{w}) \approx v_{\pi}(s)$$

Linear value function approx.

$$\hat{v}(s, \mathbf{w}) \doteq \sum w_i x_i(s) \\ = \langle \mathbf{w}, \mathbf{x}(s) \rangle$$

Nonlinear value function approx. w/ NN

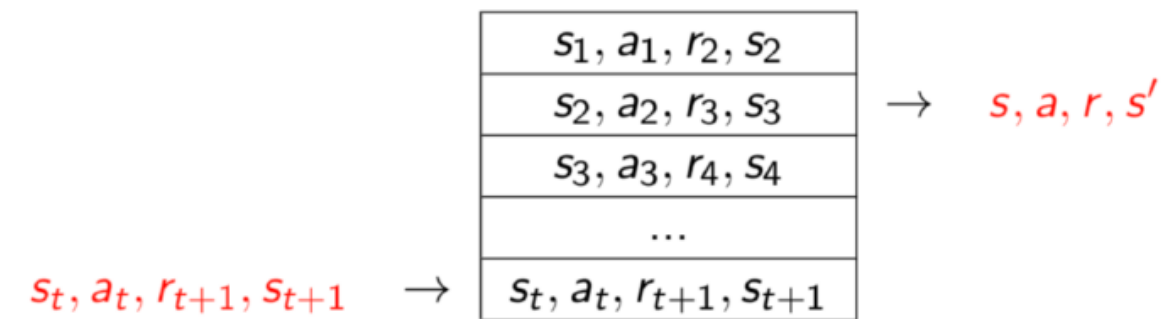


-change of data distribution may lead to oscillation and uncertainty during training and result in local minimum convergence

Deep Q-Networks (DQN)

Deep Q-Networks (DQN)

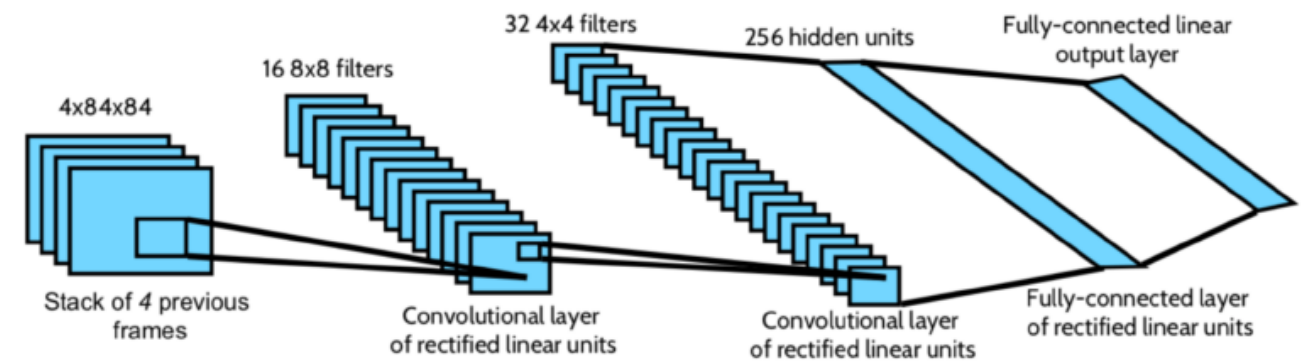
Replay buffer



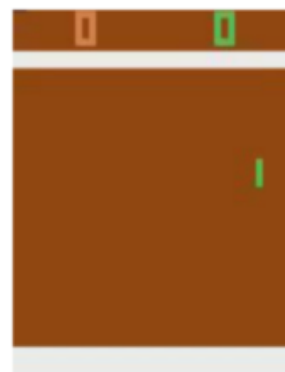
To remove correlations btw samples.

Update with fixed target network θ^-

$$l = \left(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2$$



DQN on Atari (Network Architecture)



Pong



Enduro



Beamrider



Q*bert

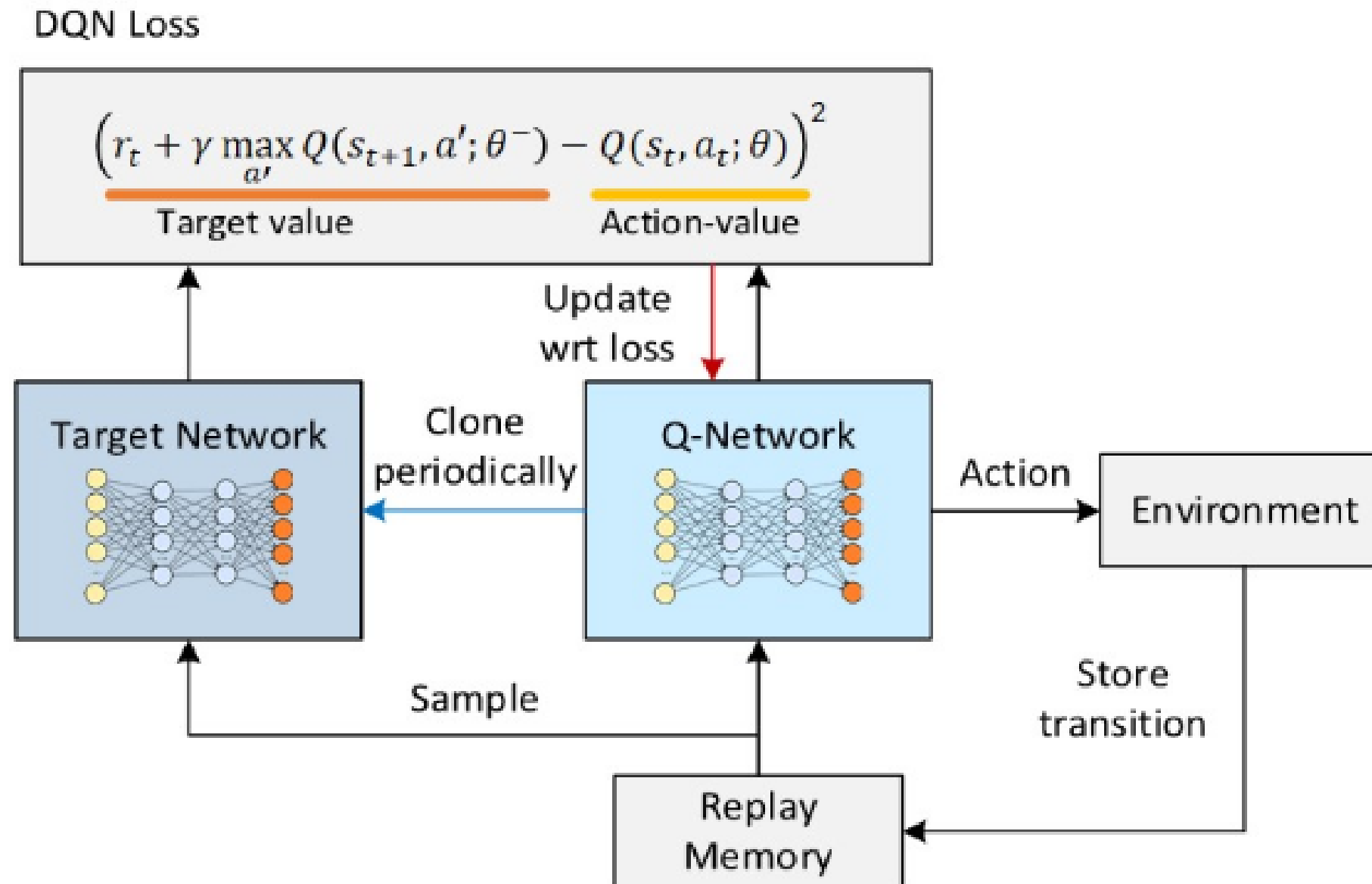
Atari games

-experience replay and target network resolves the problem of Deep Q-learning

Deep Q-Networks (DQN)

Deep Q-Networks (DQN)

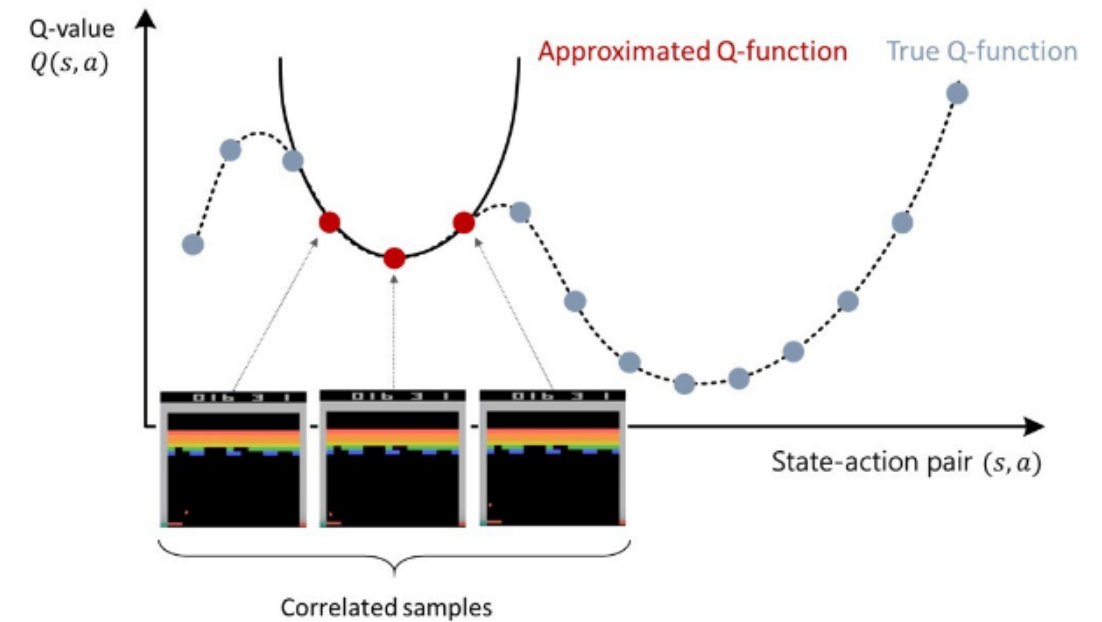
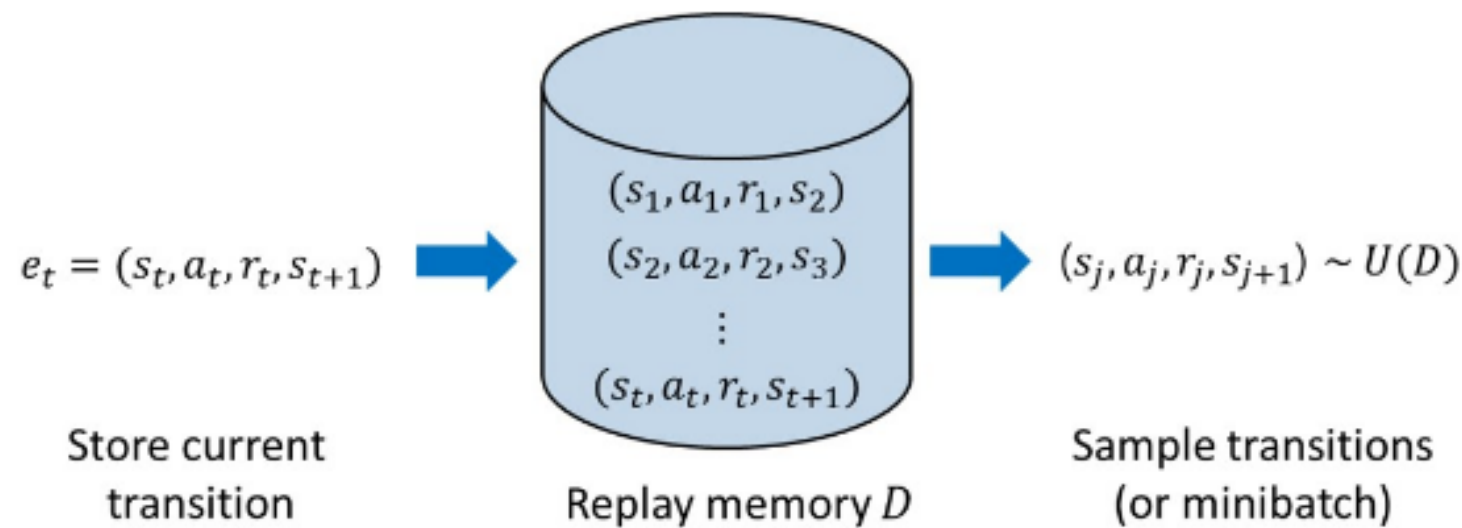
- Target network



Deep Q-Networks (DQN)

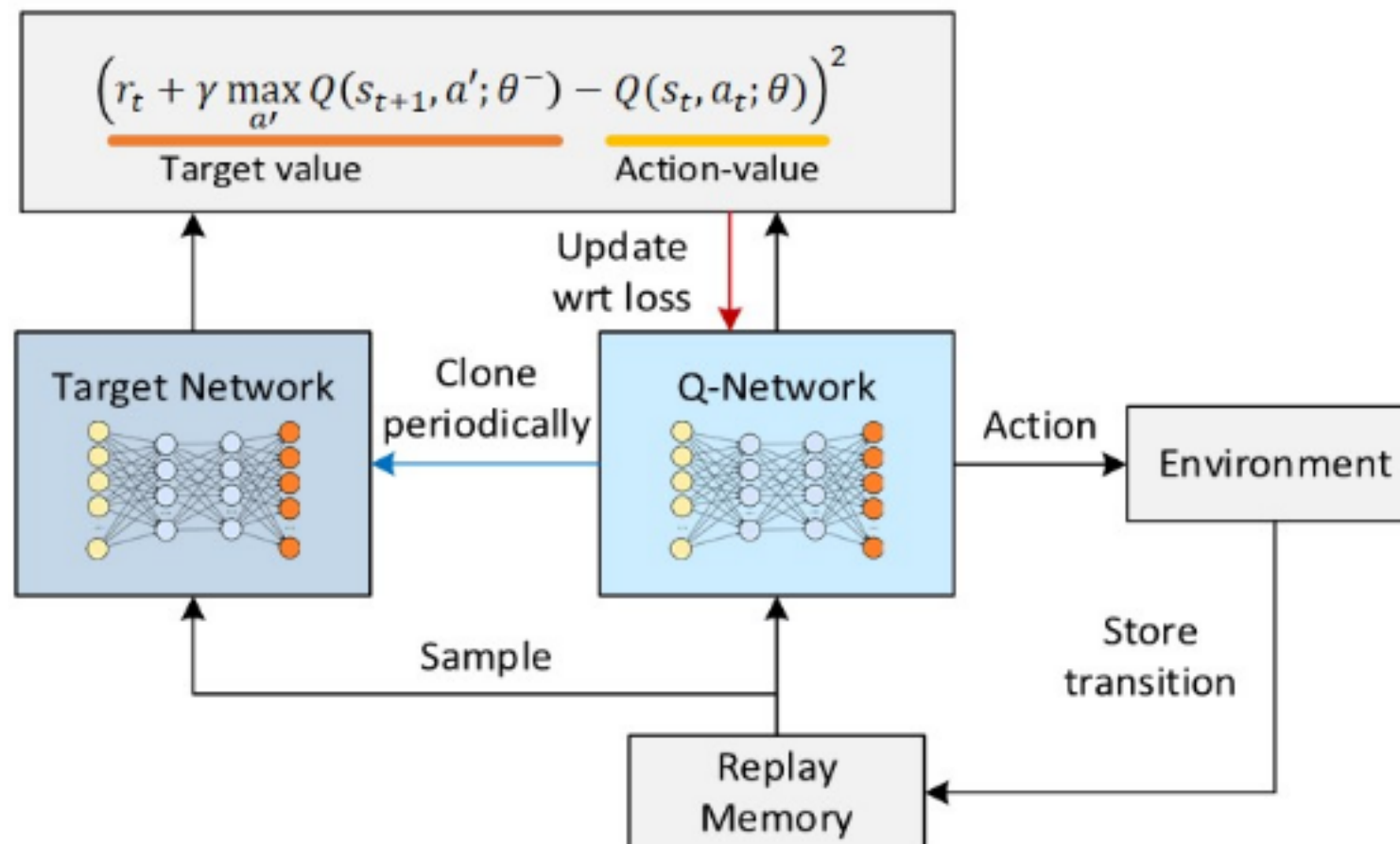
Deep Q-Networks (DQN)

- Experience replay



- Target Network

DQN Loss



Policy Gradient Theorem

Policy Gradient Theorem

Policy

$$\pi(a | s) = \text{Pr}(A_t = a | S_t = s)$$

↓
 π_θ

Reward function (policy objective function)

$$J(\theta) = \sum_{s \in S} d^\pi(s) V^\pi(s) = \sum_{s \in S} d^\pi(s) \sum_{a \in A} \pi_\theta(a | s) Q^\pi(s, a)$$

Stationary distribution (d^π)

$$d^\pi(s) = \lim_{t \rightarrow \infty} P(s_t = s | s_0, \pi_\theta)$$

Starting from s_0 , a probability that the state becomes s .
(under the policy π_θ)

