

Reinforce Learning

Deep Q-Learning (DQN)

소프트웨어 공대 강의

노기섭 교수

(kafa46@hongik.ac.kr)

Tackle the limitation of Dynamic Programming

Solution

- Create an **approximator** that simulate $Q(s, a)$
 - ✓ Q-learning
- We can build the approximator using DNN \Rightarrow Deep Q Network

DNN params (weights) Bellman Equation

$$Q(s, a ; \theta) \approx Q^*(s, a)$$

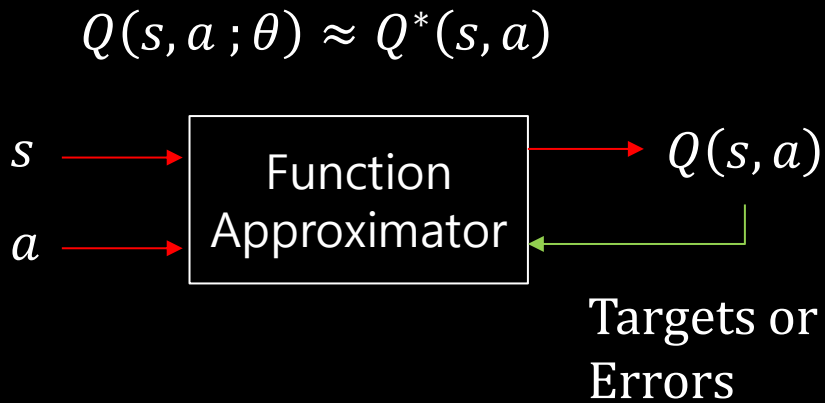
Output from DNN
(approximator)

Ground Truth

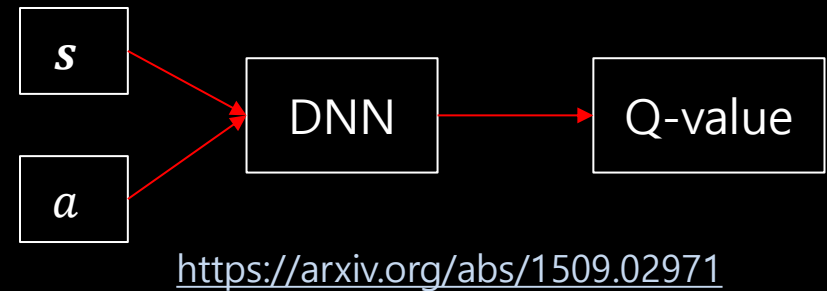
Let the DNN learn Q functions
that follow Bellman's equations.

How to Implement DQN?

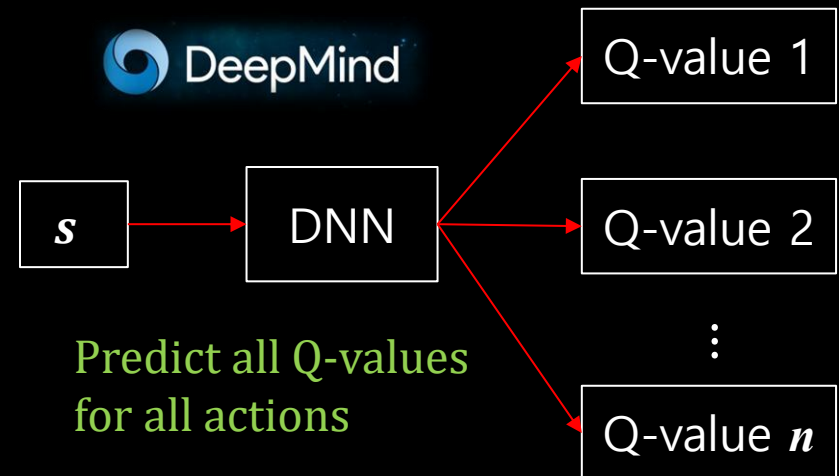
2 Solutions are possible!



Solution 1

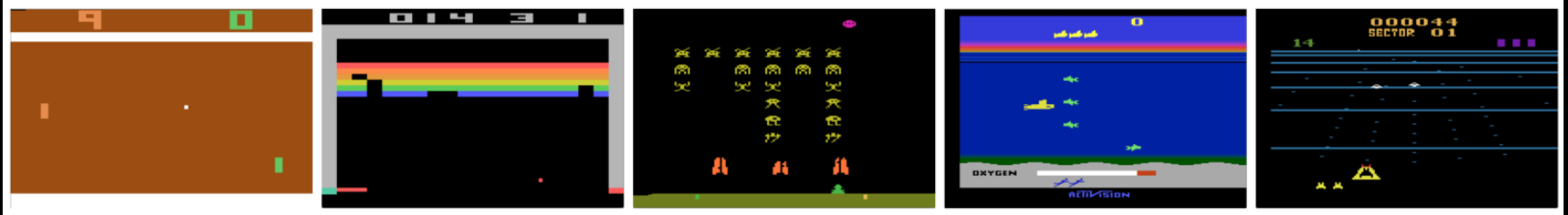


Solution 2



DeepMind DQN: playing Atari games

Atari 2600 Games



Pong

Breakout

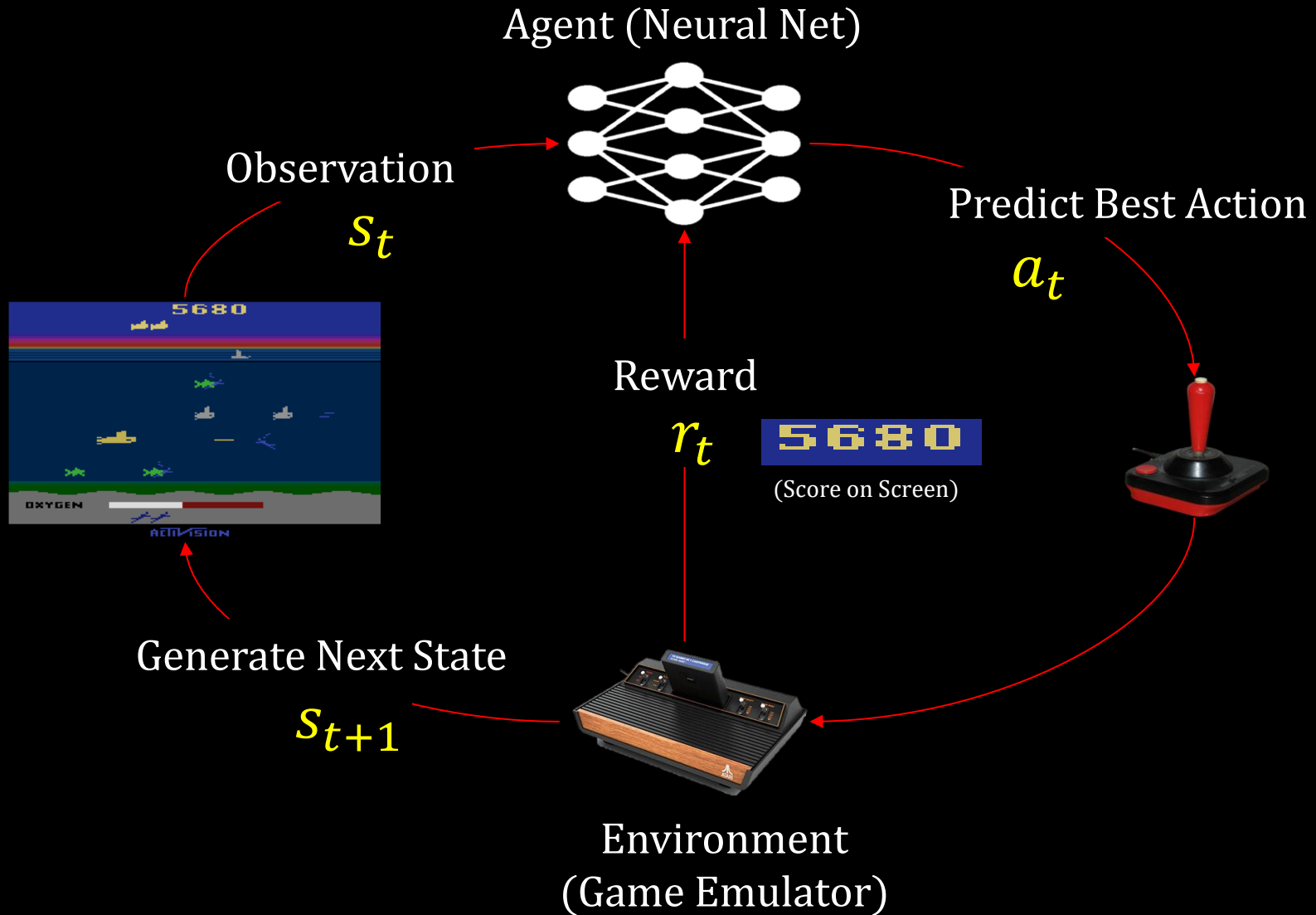
Space Invaders

Seaquest

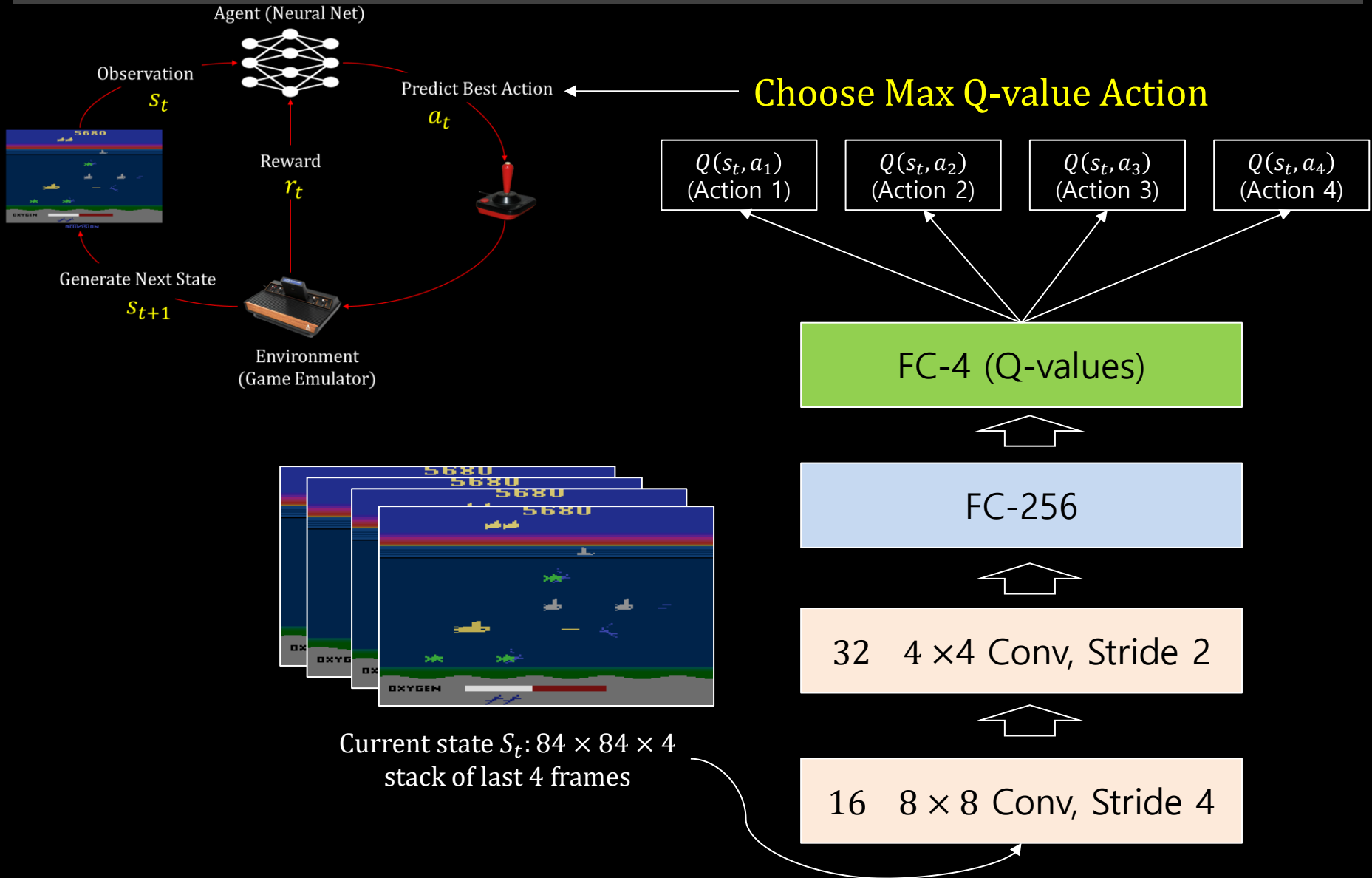
Beam Rider

- **Objective:** End game with Max Score
- **State:** Raw pixel distribution of screen
- **Action:** Game control (left, right, up, down)
- **Reward:** Score value on screen for each time step

DeepMind DQN: Algorithm Schetch



Detailed DQN Architecture



Experience Replay Buffer

- Issue: **strong correlation** between samples

$s_0 a_1 r_0 s_1 \quad s_1 a_1 r_1 s_2 \quad \dots \quad s_{t-1} a_{t-1} r_{t-1} s_t$



$(s_0 a_1 r_0 s_1), (s_1 a_1 r_1 s_2), (s_2 a_2 r_2 s_3), \dots, (s_{t-1} a_{t-1} r_{t-1} s_t)$

- ✓ Current params determine the next training samples
- ✓ Training samples that move the robot arm to the left and then sample can only be seen around the left



Biased sample → Broken i.i.d

Hidden Problems - Tricks in DQN (1/2) i.i.d ?

Independent & Identically Distributed (i.i.d)

Independent (독립)

Each data sample does **not affect** each other.

[동전 던지기] 앞면이 나왔다고 다음에 앞면이 나올 확률이 바뀌지 않음



$$P(X_1, X_2, \dots, X_n)$$

$$= \prod_{i=1}^n P(X_i)$$

Identically Distributed (동일분포)

All data came from the same probability distribution

All samples are collected from the same sensor, same conditions, same environment, etc.

$$X_1 \sim X_2 \sim \dots \sim X_n$$

$$X_1, X_2, \dots, X_n \sim \mathcal{D}$$

$$X_i \sim \mathcal{D}$$

$\mu, \sigma, \text{distribution}$
are all same!



Why is i.i.d. important in deep learning?

SGD has only one assumption!!!

The gradients calculated from the mini-batch are "unbiased estimators" of the true gradients

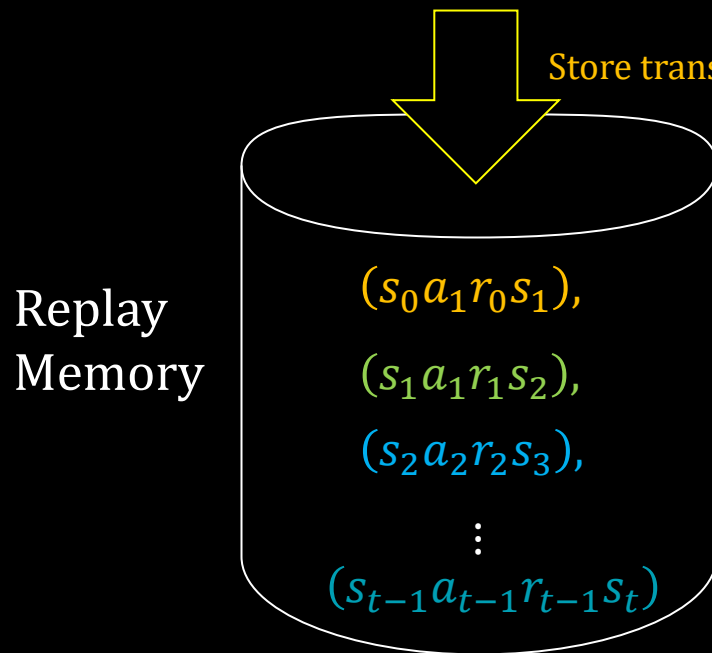
	Not Independent	Not Identically Distributed
Cause	<ul style="list-style-type: none">• Samples look like each other• Gradient over-biased in certain directions• Increased variance, wobbly orientation	<ul style="list-style-type: none">• Domain shift• Federated Learning• Changed Process Environment
Effects	<ul style="list-style-type: none">• Oscillation (진동)• Learning rate sensitivity• Unconverged loss	<ul style="list-style-type: none">• Different optimal point for each batch• Fail to converge to one minimum

결국에는 중재자 없는 정치 싸움판으로 변질 것임

Experience Replay Buffer

- Continuously save transition to replay memory while episode play

Transitions: s_t, a_t, r_t, s_{t+1}



- ✓ When Replay memory is full enough,
 - Make a random minibatch of transition
 - Sample it there to learn Q-network

[Ex.] $\left\{ \begin{array}{l} (s_7 a_7 r_7 s_8), \\ (s_{51} a_{51} r_{51} s_{52}), \\ (s_2 a_2 r_2 s_3), \\ \vdots \end{array} \right.$

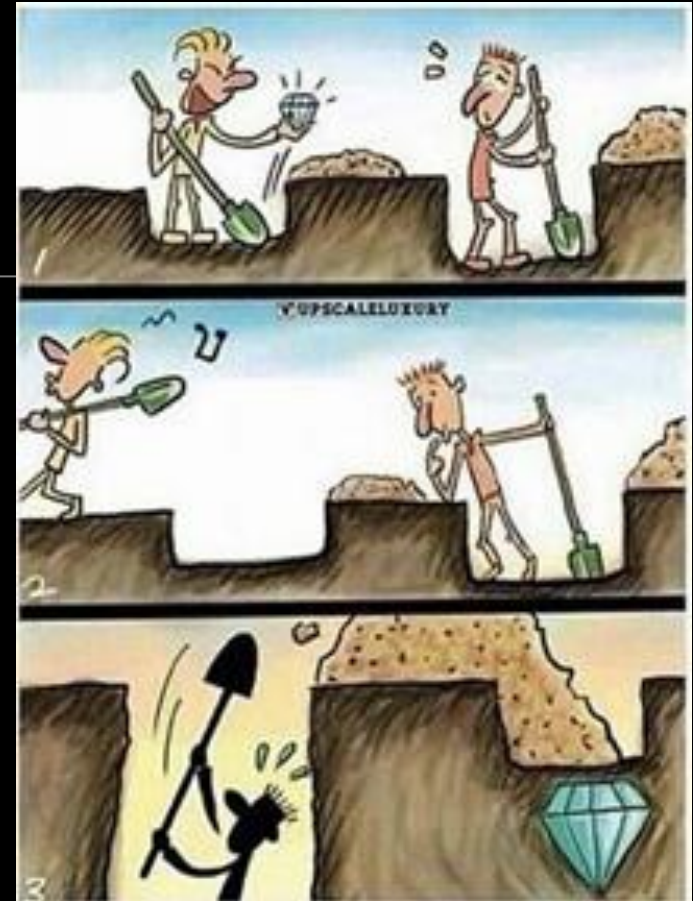
Exploitation vs. Exploration: ϵ - greedy Policy

Exploitation (활용)	Exploration (탐험)
Perform <ul style="list-style-type: none">optimal actiongreedy action	Random action

ϵ - greedy Policy

$$a_t = \begin{cases} \text{random action with prob } \epsilon \\ \arg \max_a Q(s_t, a) \text{ with prob } 1 - \epsilon \end{cases}$$

- Early: $\epsilon \approx 1.0$
(almost exploration)
- Slowly decrease ϵ
(exploration \rightarrow exploitation)
- End: $\epsilon \approx 0.0$ (almost exploitation)



<https://ai-ml-analytics.com/reinforcement-learning-exploration-vs-exploitation-tradeoff/>

Fixed Target Network

More ingredient

The calculation of the error contains a target function that varies according to network params.

- ✓ Unstable error
- ✓ Not converge



A dragon chases its own tail ^^

Target (Bellman Eq.)

Prediction

$$Q(s_t, a) \leftarrow Q(s_t, a) + \alpha \left[\underbrace{r_{t+1} + \gamma \max_p Q(s_{t+1}, p)}_{\text{Ground Truth}} - \underbrace{Q(s_t, a)}_{\text{Prediction}} \right]$$

DQN Trick: update target every n (for example, 1,000) steps!

DQN Algorithm in pseudo code

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights

for episode=1 to M do

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

for $t = 1$ to T do

With probability ϵ select random action $a_t \leftarrow$ **exploration**

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

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

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

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

Training Procedure

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

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

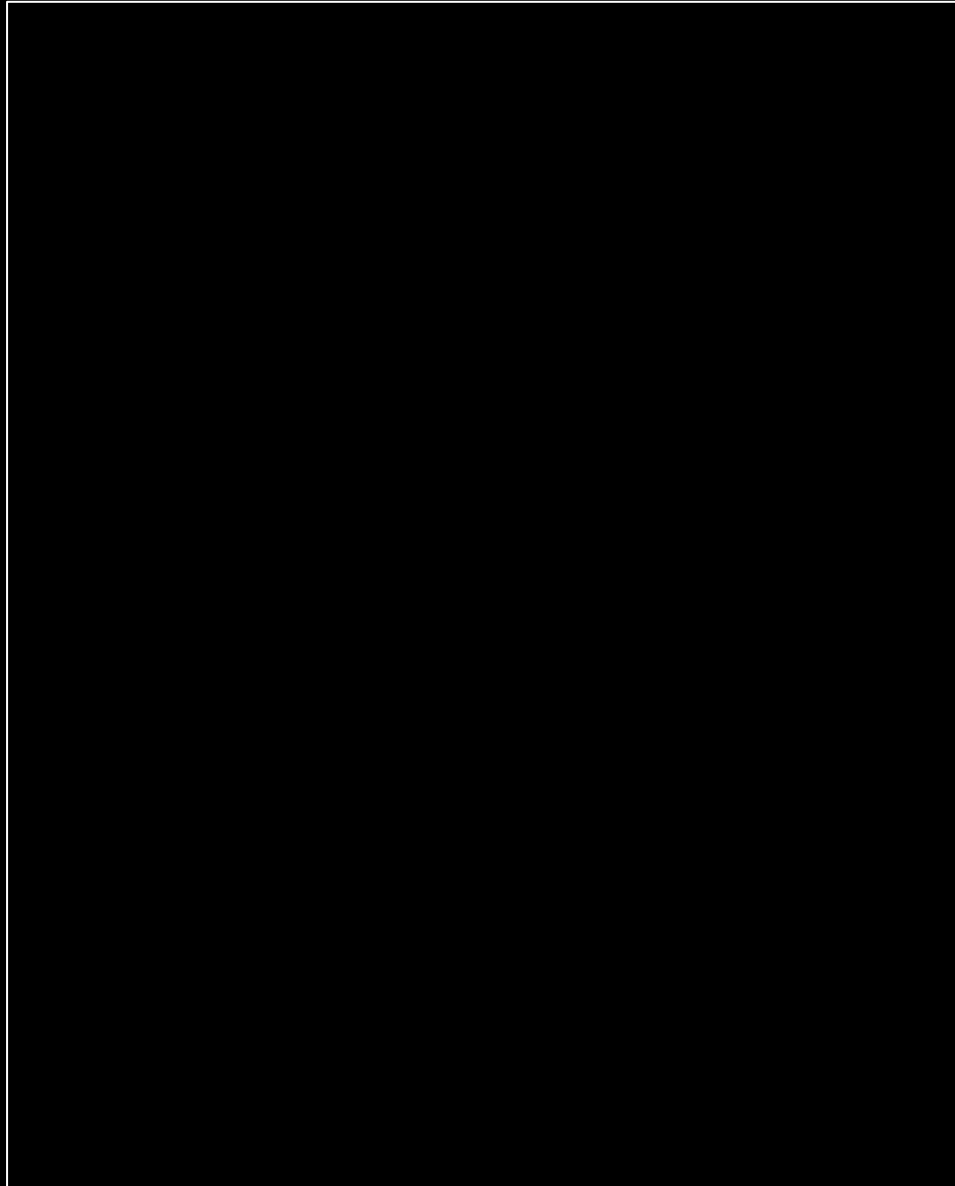
Perform a gradient descent step on $(y_i - Q(\phi_j, a_j; \theta))^2$ according to Eq. 3

Target (Bellman)

Prediction (output) from NN

Google DeepMind DQN - Atari Breakout

Applying
DQN into
Atari Game



영상 출처:
[https://youtu.be/
V1eYniJ0Rnk?si=
WdQARJq09UYM](https://youtu.be/V1eYniJ0Rnk?si=WdQARJq09UYM)

Summary

■ Q-learning

- Idea: Learn a function that approximates $Q(s, a)$ using the Bellman equation.
- Pros: Highly sample-efficient compared to policy optimization methods due to effective data reuse.
- Cons: Optimizes the policy indirectly, which can be less stable and less reliable.
- The learned Q-function is only an approximation and cannot be guaranteed to be perfect.
- Main challenge: Exploration.

■ DQN (Deep Q-Network)

■ Next Step: Policy Gradient

- Directly solves an optimization problem and converges to a local minimum of the cost function.
- Requires more samples than DQN.



수고하셨습니다 ..^^..