# Dueling Network Architectures for Deep Reinforcement Learning (ICML 2016)

Yoonho Lee

Department of Computer Science and Engineering
Pohang University of Science and Technology

October 11, 2016

*POSTECH*

# Outline

*POSTECH*

# Outline

*POSTECH*

# Definition of RL

general setting of RL



observation $O_t$

action $A_t$

reward $R_t$

POSTECH

# Definition of RL

atari setting



observation
$O_t$

action
$A_t$

reward  $R_t$

# Outline

POSTECH

# Mathematical formulations

### Definition

Return $G_t$ is the cumulative discounted reward from time t

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \ldots$$

POSTECH

# Mathematical formulations

### Definition
Return $G_t$ is the cumulative discounted reward from time t

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \ldots$$

### Definition
A policy $\pi$ is a function that selects actions given states

$$\pi(s) = a$$

- ▶ The goal of RL is to find $\pi$ that maximizes $G_0$

*POSTECH*

# Mathematical formulations
## Q-Value

$$G_t = \sum_{i=0}^{\infty} \gamma^k r_{t+i+1}$$

### Definition
The action-value (Q-value) function $Q_\pi(s, a)$ is the expectation of $G_t$ under taking action a, and then following policy $\pi$ afterwards

$$Q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a, A_{t+i} = \pi(S_{t+i}) \forall i \in \mathbb{N}]$$

► "How good is action a in state s?"

# Mathematical formulations

### Definition

The optimal Q-value function $Q_*(s, a)$ is the maximum Q-value over all policies

$$Q_*(s, a) = \max_\pi Q_\pi(s, a)$$

# Mathematical formulations

### Definition
The optimal Q-value function $Q_*(s, a)$ is the maximum Q-value over all policies

$$Q_*(s, a) = \max_\pi Q_\pi(s, a)$$

### Theorem
There exists a policy $\pi_*$ such that $Q_{\pi_*}(s, a) = Q_*(s, a)$ for all $s, a$

- ▶ Thus, it suffices to find $Q_*$

POSTECH

# Mathematical formulations

### Bellman Optimality Equation

$Q_*(s, a)$ satisfies the following equation:

$$Q_*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_*(s', a')$$

POSTECH

# Mathematical formulations

## Bellman Optimality Equation

$Q_*(s, a)$ satisfies the following equation:

$$Q_*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_*(s', a')$$

## Q-Learning

Let $a$ be $\epsilon$-greedy w.r.t. $Q$, and $a'$ be optimal w.r.t. $Q$. $Q$ converges to $Q_*$ if we iteratively apply the following update:

$$Q(s, a) \leftarrow \alpha(R(s, a) + \gamma Q(s', a')) + (1 - \alpha)Q(s, a)$$

POSTECH

# Mathematical formulations

### Value-Based RL

- Estimate $Q_*(s, a)$
- Deep Q Network

### Policy-Based RL

- Search directly for optimal policy $\pi_*$
- DDPG, TRPO...

### Model-Based RL

- Use the (learned or given) transition model of environment
- Tree Search, DYNA ...

*POSTECH*

# Outline

POSTECH

# Neural Fitted $Q$ Iteration

---

**NFQ_main()** {
input: a set of transition samples $D$; output: Q-value function $Q_N$
    k=0
    init_MLP() $\to Q_0$;
    Do {
        generate_pattern_set $P = \{(input^l, target^l), l = 1, \ldots, \#D\}$ where:
            $input^l = s^l, u^l,$
            $target^l = c(s^l, u^l, s'^l) + \gamma\, min_b Q_k(s'^l, b)$
        Rprop_training($P$) $\to Q_{k+1}$
        k:= k+1
    } WHILE $(k < N)$

---

**Fig. 1.** Main loop of NFQ

▶ Supervised learning on (s,a,r,s')

POSTECH

# Neural Fitted $Q$ Iteration

input

---

**NFQ_main()** {
input: a set of transition samples $D$; output: Q-value function $Q_N$
    k=0
    init_MLP() $\rightarrow Q_0$;
    Do {
        generate_pattern_set $P = \{(input^l, target^l), l = 1, \ldots, \#D\}$ where:
            $input^l = s^l, u^l,$
            $target^l = c(s^l, u^l, s'^l) + \gamma \min_b Q_k(s'^l, b)$
        Rprop_training($P$) $\rightarrow Q_{k+1}$
        k:= k+1
    } WHILE $(k < N)$

---

**Fig. 1.** Main loop of NFQ

# Neural Fitted $Q$ Iteration

target equation

---

**NFQ_main()** {
input: a set of transition samples $D$; output: Q-value function $Q_N$
    k=0
    init_MLP() $\rightarrow Q_0$;
    Do {
        generate_pattern_set $P = \{(input^l, target^l), l = 1, \ldots, \#D\}$ where:
            $input^l = s^l, u^l,$
            $target^l = c(s^l, u^l, s'^l) + \gamma \min_b Q_k(s'^l, b)$
        Rprop_training($P$) $\rightarrow Q_{k+1}$
        k:= k+1
    } WHILE $(k < N)$

---

**Fig. 1.** Main loop of NFQ

POSTECH

# Neural Fitted $Q$ Iteration

## Shortcomings

- Exploration is independent of experience
- Exploitation does not occur at all
- Policy evaluation does not occur at all
- Even in exact(table lookup) case, not guaranteed to converge to $Q_*$

*POSTECH*

# Outline

POSTECH

# Deep $Q$ Network

## action policy

**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 $\varepsilon$ select a random action $a_t$
        otherwise select $a_t = \arg\max_a Q(\phi(s_t), a; \theta)$
        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 $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$
        Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$
        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$
        Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the network parameters $\theta$
        Every $C$ steps reset $\hat{Q} = Q$
    **End For**
**End For**

▶ Choose an $\epsilon$-greedy policy w.r.t. $Q$

# Deep $Q$ Network

## network freezing

**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** $t = 1, \text{T}$ **do**
        With probability $\varepsilon$ 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 $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$
        Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$
        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} Q\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$
        Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the
        network parameters $\theta$
        Every $C$ steps reset $\hat{Q} = Q$
    **End For**
**End For**

▶ Get a copy of the network every $C$ steps for stability

# Deep $Q$ Network

**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 $\varepsilon$ 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 $\left(\phi_t,a_t,r_t,\phi_{t+1}\right)$ in $D$
  Sample random minibatch of transitions $\left(\phi_j,a_j,r_j,\phi_{j+1}\right)$ from $D$
  Set $y_j = \begin{cases} r_j & \text{if episode terminates at step j+1} \\ r_j + \gamma \ \text{max}_{a'} \ \hat{Q}\left(\phi_{j+1},a'; \theta^-\right) & \text{otherwise} \end{cases}$
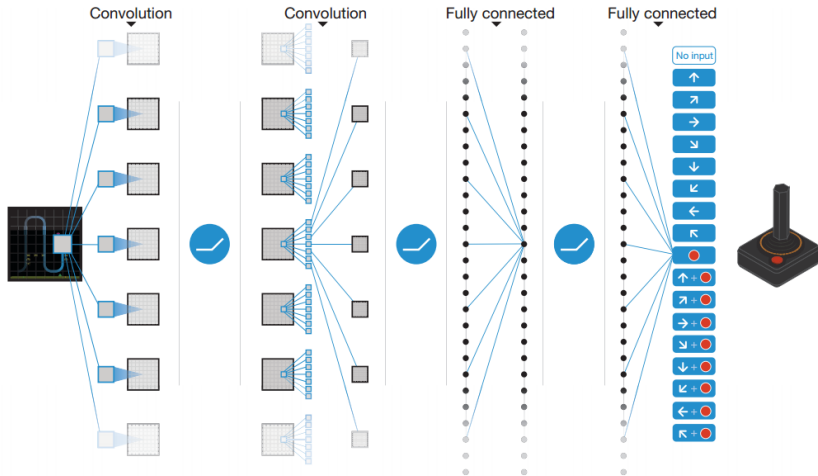  Perform a gradient descent step on $\left(y_j - Q\left(\phi_j,a_j; \theta\right)\right)^2$ with respect to the
  network parameters $\theta$
  Every $C$ steps reset $\hat{Q} = Q$
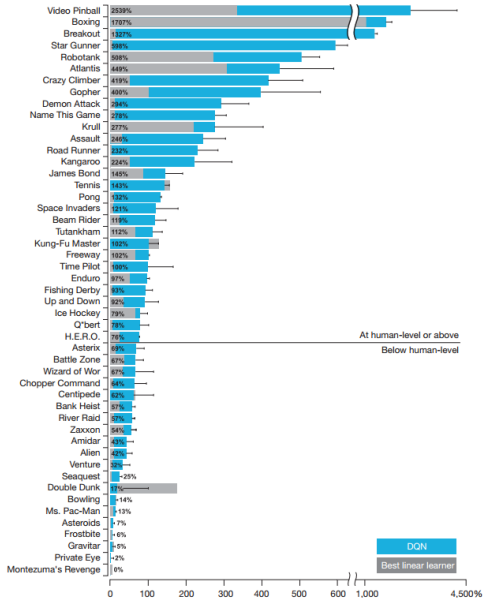 **End For**
**End For**

# Deep $Q$ Network

# Deep Q Network

## performance

# Deep $Q$ Network

overoptimism

**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 $\varepsilon$ 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 $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$
        Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$
        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} Q\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$
        Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the
        network parameters $\theta$
        Every $C$ steps reset $\hat{Q} = Q$
    **End For**
**End For**

▶ What happens when we overestimate $Q$?

# Deep $Q$ Network

### Problems

- Overestimation of the $Q$ function at any $s$ spills over to actions that lead to $s \rightarrow$ Double DQN
- Sampling transitions uniformly from $D$ is inefficient $\rightarrow$ Prioritized Experience Replay

# Outline

POSTECH

# Double Deep $Q$ Network

target

We can write the DQN target as:

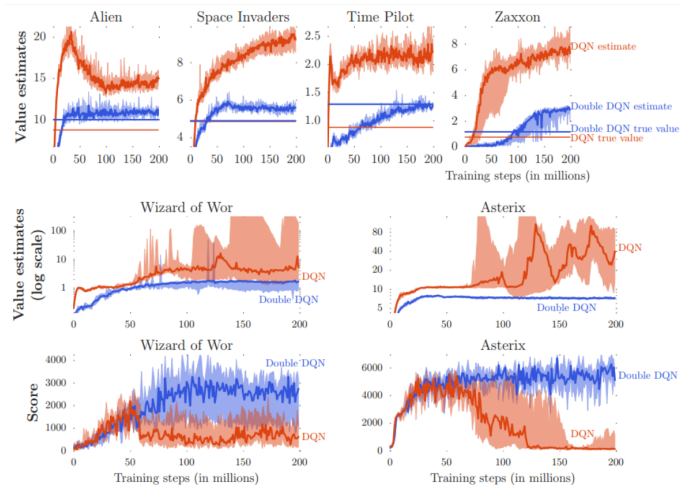$$Y_t^{DQN} = R_{t+1} + \gamma Q(S_{t+1}, \underset{a}{argmax}\, Q(S_{t+1}, a; \theta_t^-); \theta_t^-)$$

Double DQN's target is:

$$Y_t^{DDQN} = R_{t+1} + \gamma Q(S_{t+1}, \underset{a}{argmax}\, Q(S_{t+1}, a; \theta_t); \theta_t^-)$$

This has the effect of decoupling action selection and action evaluation

# Double Deep $Q$ Network

performance



- Much stabler with very little change in code

POSTECH

# Outline

*POSTECH*

# Prioritized Experience Replay

**Algorithm 1** Double DQN with proportional prioritization

1: **Input:** minibatch $k$, step-size $\eta$, replay period $K$ and size $N$, exponents $\alpha$ and $\beta$, budget $T$.
2: Initialize replay memory $\mathcal{H} = \emptyset$, $\Delta = 0$, $p_1 = 1$
3: Observe $S_0$ and choose $A_0 \sim \pi_\theta(S_0)$
4: **for** $t = 1$ **to** $T$ **do**
5:     Observe $S_t, R_t, \gamma_t$
6:     Store transition $(S_{t-1}, A_{t-1}, R_t, \gamma_t, S_t)$ in $\mathcal{H}$ with maximal priority $p_t = \max_{i<t} p_i$
7:     **if** $t \equiv 0 \mod K$ **then**
8:         **for** $j = 1$ **to** $k$ **do**
9:             Sample transition $j \sim P(j) = p_j^\alpha / \sum_i p_i^\alpha$
10:            Compute importance-sampling weight $w_j = (N \cdot P(j))^{-\beta} / \max_i w_i$
11:            Compute TD-error $\delta_j = R_j + \gamma_j Q_{\text{target}}(S_j, \arg\max_a Q(S_j, a)) - Q(S_{j-1}, A_{j-1})$
12:            Update transition priority $p_j \leftarrow |\delta_j|$
13:            Accumulate weight-change $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_\theta Q(S_{j-1}, A_{j-1})$
14:         **end for**
15:         Update weights $\theta \leftarrow \theta + \eta \cdot \Delta$, reset $\Delta = 0$
16:         From time to time copy weights into target network $\theta_{\text{target}} \leftarrow \theta$
17:     **end if**
18:     Choose action $A_t \sim \pi_\theta(S_t)$
19: **end for**

▶ Update 'more surprising' experiences more frequently

# Outline

**POSTECH**

# Dueling Network



▶ The scalar approximates $V$, and the vector approximates $A$
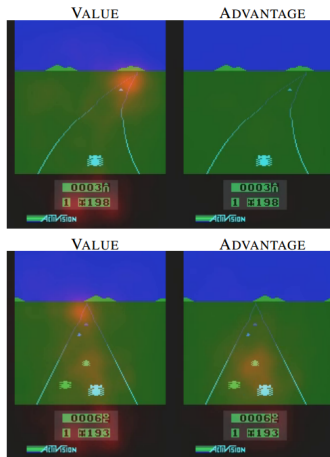
POSTECH

# Dueling Network

The exact forward pass equation is:

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + (A(s, a; \theta, \alpha) - \max_{a'} A(s, a'; \theta, \alpha))$$

The following module was found to be more stable without losing much performance:

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + (A(s, a; \theta, \alpha) - \frac{1}{|A|} \sum_{a'} A(s, a'; \theta, \alpha))$$

POSTECH

# Dueling Network

- Advantage network only pays attention when current action crucial

# Dueling Network

performance

| | 30 no-ops | | Human Starts | |
|---|---|---|---|---|
| | **Mean** | **Median** | **Mean** | **Median** |
| Prior. Duel Clip | **591.9%** | **172.1%** | **567.0%** | **115.3%** |
| Prior. Single | 434.6% | 123.7% | 386.7% | 112.9% |
| Duel Clip | **373.1%** | **151.5%** | **343.8%** | **117.1%** |
| Single Clip | 341.2% | 132.6% | 302.8% | 114.1% |
| Single | 307.3% | 117.8% | 332.9% | 110.9% |
| Nature DQN | 227.9% | 79.1% | 219.6% | 68.5% |

▶ Achieves state of the art in the Atari domain among DQN algorithms

POSTECH

# Dueling Network

## Summary

- Since this is an improvement only in network architecture, methods that improve DQN(e.g. Double DQN) are all applicable here as well
- Solves problem of $V$ and $A$ typically being of different scale
- Updates $Q$ values more frequently than a single-stream DQN, where only a single $Q$ value is updated for each observation
- Implicitly splits the credit assignment problem into a recursive binary problem of "now or later"

*POSTECH*

# Dueling Network

## Shortcomings

- Only works for $|A| < \infty$
- Still not able to solve tasks involving long-term planning
- Better than DQN, but sample complexity is still high
- $\epsilon$-greedy exploration is essentially random guessing