

Double Deep Q-Network Solution

Yermek Kapushev

Skoltech

2017

Q-learning

- ▶ Under a given policy π , the true value of action a in a state s is

$$Q_\pi(s, a) = \mathbb{E}[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | s_0 = s, a_0 = a, \pi],$$

where $\gamma \in [0, 1]$.

- ▶ The optimal value is

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

- ▶ An optimal policy is derived from $Q^*(s, a)$ by taking
$$\arg \max_a Q^*(s, a)$$

Learning optimal action-value function

Let $Q(s, a, \theta)$ be a model, that approximates Q-function. θ is a vector of its parameters.

1. Initialize $Q(s, a; \theta)$

2. Observe:

$$s' \leftarrow s, \quad a' \leftarrow a$$

get new state s , get reward r .

3. Update weights

$$\theta \leftarrow \theta + \alpha(r + \gamma \max_{\tilde{a}} Q(s, \tilde{a}, \theta) - Q(s', a', \theta)) \nabla Q(s', a', \theta)$$

4. Go to step 2.

RL is known to be unstable, or even to diverge when a nonlinear function approximator such as an ANN is used to represent Q-function:

- ▶ Correlations in the sequence of observations
- ▶ Small updates of Q may significantly change the policy \Rightarrow change of data distribution
- ▶ Correlations between action values $Q(s, a)$ and target values
$$r + \gamma \max_a Q(s, a)$$

Solution:

- ▶ Randomization over the data *experience replay* \Rightarrow remove correlations in the observations sequence and smoothing over changes in the data distribution.
- ▶ *Iterative update* of action-values and target values \Rightarrow reduce correlations between them

Experience replay

- ▶ Store the agents experience $e_t = (s_t, a_t, r_t, s_{t+1})$ at each time step in a data set $D = \{e_1, \dots, e_t\}$
- ▶ Apply Q-learning updates on mini-batches of experience $(s, a, r, s') \sim \text{Uniform}(D)$

Double Q-learning

- ▶ The following loss function is used:

$$L(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a', \theta_i^-) - Q(s, a, \theta_i) \right)^2 \right],$$

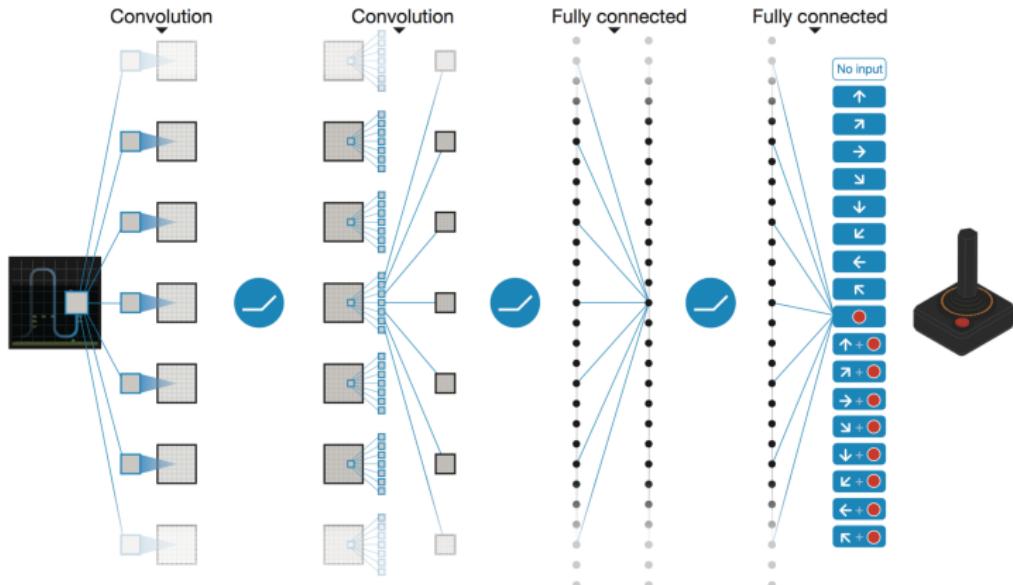
where θ_i^- — parameters used to compute action-value, θ_i — parameters used to compute target value.

- ▶ Parameters θ_i^- are updated every C steps with the parameters θ_i .

Learning to play Atari games

- ▶ Arcade Learning Environment (ALE) — python API to ROMs emulator.
- ▶ ANN is used for Double DQN
- ▶ Input data: 210 * 160 video at 60 Hz
- ▶ Output data: action, game score
- ▶ Internal game state is not observed. Only frames (images) are observed.

Deep convolutional network architecture



Input: 4 frames of $84 * 84$

Output: single output for each possible action

