# ACTOR CRITIC METHODS IN SIMULATING SIMPLE NATURAL SELECTION ENVIRONMENT

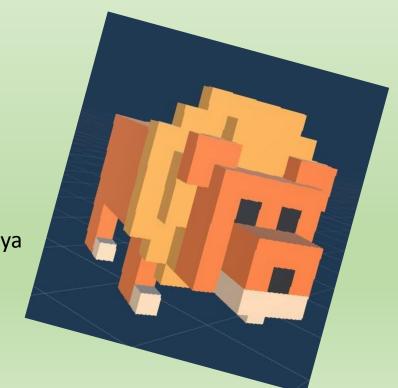




Ümit Akköse 19190

19233

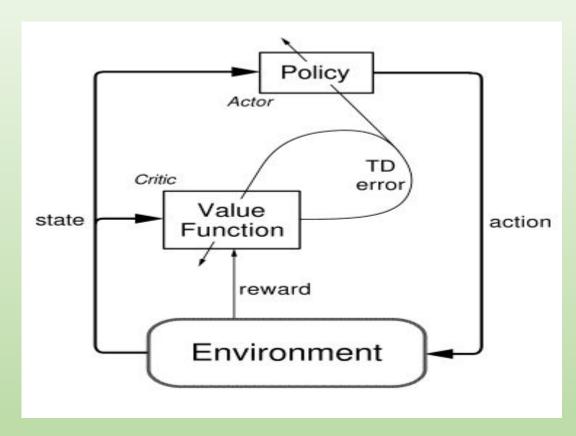
Barış Temel Veysel Oğulcan Kaya 17804



## What is Deep Reinforcement Learning?

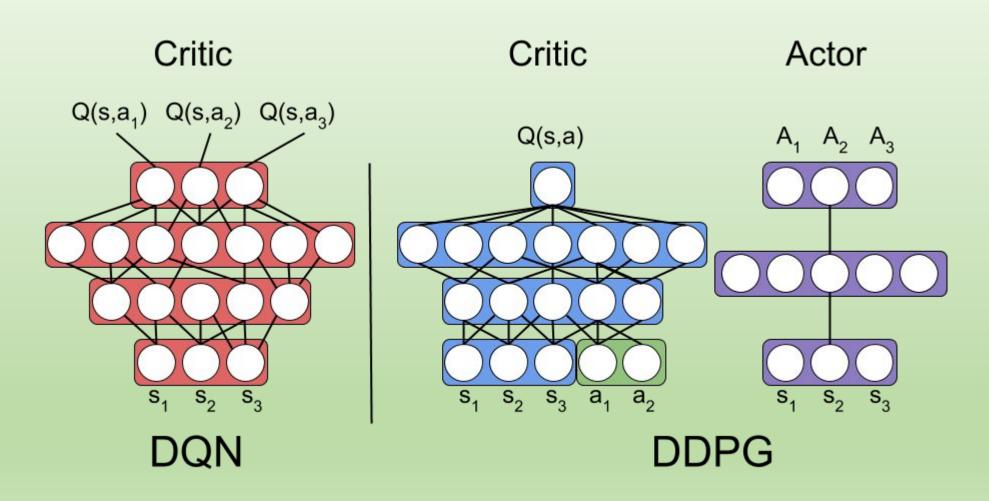
- Intelligent machines can learn from their actions similar to the way humans learn from experience.
- Inherent in this type of machine learning is that an agent is rewarded or penalised based on their actions.
- Actions that get them to the target outcome are rewarded (reinforced).

#### **Actor-Critic Method**

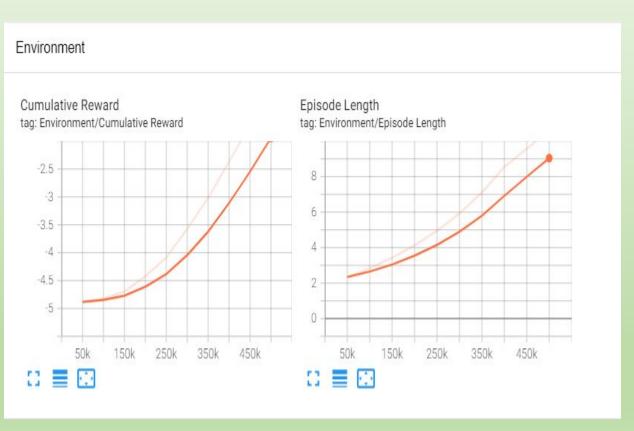


- The policy structure is known as the actor, because it is used to select actions, and the estimated value function is known as the critic.
- The critic is a state-value function.

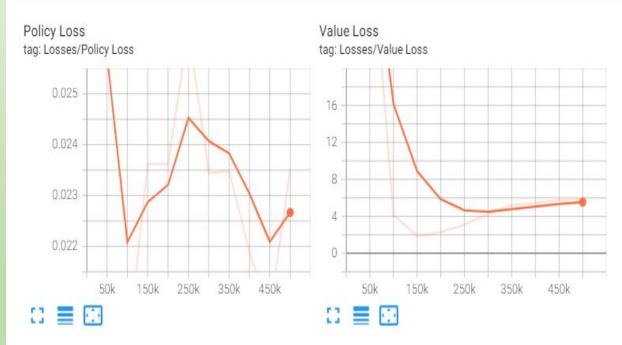
#### **DQN&DDPG STRUCTURE**



#### **DDPG** Results







### The pseudo-algorithm of DDPG

- Initialize actor network μ<sub>θ</sub> and critic Q<sub>Φ</sub> with random weights.
- Create the target networks  $\mu_{\theta}$  0 and  $Q_{\phi 0}$ .
- Initialize experience replay memory D of maximal size N.
- for episode  $\in$  [1, M]:
  - $\circ$  Initialize random process ξ.
  - Observe the initial state so.
  - o for t  $\subseteq$  [0, Tmax]:
    - Select the action at =  $\mu\theta$ (st) + ξ according to the current policy and the noise.
    - Perform the action at and observe the next state st+1 and the reward rt+1.
    - Store (st, at, rt+1, st+1) in the experience replay memory.
    - Sample a minibatch of N transitions randomly from D.
    - For each transition (sk, ak, rk, s0 k) in the minibatch: Compute the target value using target networks yk = rk + γ Qφ0(s 0 k ,  $\mu\theta$  0(s 0 k )).
    - Update the critic by minimizing:  $L(\phi) = 1 \text{ N X k (yk } Q\phi(sk, ak))2$
    - ♦ **Update** the actor using the sampled policy gradient:  $\nabla \theta J(\theta) = 1 \text{ N X k } \nabla \theta \mu \theta(\text{sk}) \times \nabla a Q \phi(\text{sk, a}) | a = \mu \theta(\text{sk})$
    - **♦ Update** the target networks:  $\theta$  0 ←  $\tau\theta$  + (1 −  $\tau$ )  $\theta$  0  $\phi$  0 ←  $\tau\phi$  + (1 −  $\tau$ )  $\phi$

