

Parameter Name:

S: current State, A: current Action, R: Reward, S'=next State

W=weights of Separate target Network

Q_Target_Next Equation And Q_Target:

$$Q_{\text{target_next}} = \max_a \hat{Q}(S', A, W)$$

$$Q_{\text{target}} = \text{reward} + (\gamma * Q_{\text{target_next}} * (1 - \text{done}))$$

My Loss function and Optimizer is

$$\text{Loss function} = (1/2) * (Q_{\text{Expected}} - Q_{\text{target}})^2$$

Optimizer: Adam, Learning rate: 0.0005

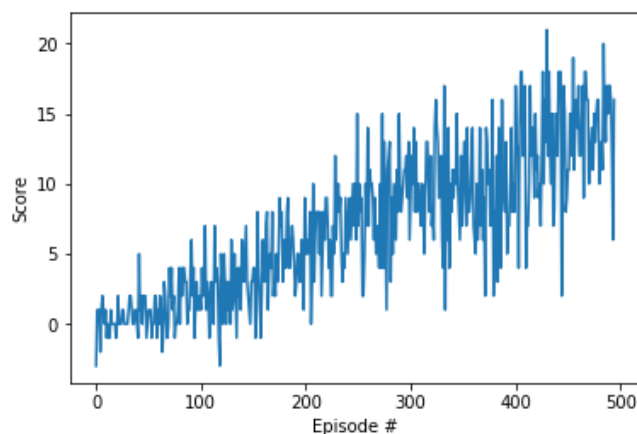
Batch_size=64:

My Qnetwork consists of this.

Input Layer vector:37 , hidden1 Layer size: 64 ,hidden2 Layer size:64, output Layer size: 4

My Plot of Rewards

Episode 100	Average Score: 0.97
Episode 200	Average Score: 3.52
Episode 300	Average Score: 7.54
Episode 400	Average Score: 9.50
Episode 495	Average Score: 13.03
Environment solved in 395 episodes!	Average Score: 13.03
Score: 16.0	



DQN Algorithms

If(memory size < BATCHSIZE)

Choose action A from state S using policy epsilon-greedy action.

Take action A, Observe reward R, input frame xt

Prepare next state: Next state $s' = \phi(x_t, x_{t+1}, x_{t+2}, x_{t+3})$

Store experience tuple (S', A, R, S) in memory D

Else

Obtain minibatch of tuples (s_j, a_j, r_j, s_{j+1}) from D

$Q_target_next = \max_a \hat{q}(S', A, W')$

$Q_target = reward + (\gamma * Q_target_next * (1 - done))$

Update $\nabla w = \alpha * (y_j - \hat{q}(S', A, W')) \nabla \hat{q}(S', A, W')$

$W' = W$

Performance Improve

At First, I want to try Deep layer and big size. And Dqn is weakness that is overestimated Q-value. Because choosing max q-value. If I use Dueling DQN, I will get more good performance than dqn. Because Dueling Dqn is separated value function and advantage function.