# Comparing deep reinforcement algorithm on lunar space lander.

1
2  **Sai Charith Reddy Pasula**                                **Nikhil Tej Kantu**
3  50320452                                                           50336812
4  *saichari@buffalo.edu*                                   *nikhilte@buffalo.edu*
5  *Department of Computer Science*                *Department of Robotics*
6  *University at buffalo*                                      *University at buffalo*

7    ## Abstract

8    The aim of the assignment is to compare the different deep reinforcement algorithm
9    on the lunar lander environment and Acrobot environment provided by openAI gym.

10

11  # 1   Introduction

12

13  ## 1.1   DQN

14
15  we approximate the Q- values by adding a term w to it.
16

$$Q(s, a, w) \approx Q^{\pi}(s, a)$$

17
18                              *Q value approximation*

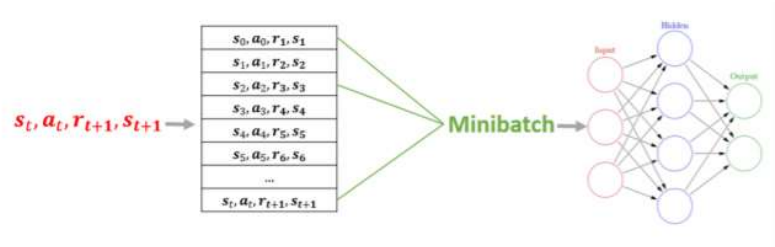19  We minimize w by using the objective function:

$$\mathcal{L}(w) = \mathbb{E}\left[\left(\underbrace{r + \gamma \max_{a'} Q(s', a', w)}_{target} - Q(s, a, w)\right)^2\right]$$

20
21                              *Objective function*
22

23  We optimise this objective function by using stochastic gradient descent.
24  Deepmind team gave three different optimizations to the original Deep RL
25  techniques:
26  • The data in the RL environment is sequential in nature. This is because the
27    agent performs one action and the next action is related to the state he ends up
28    in. So, the input to the network is correlated. To overcome this problem, we use
29    experience replay. In experience replay, the agent performs an action  according
30    to the epsilon-greedy policy. We then store the state, action, reward, next_state
31    tuple in replay memory. we sample a minibatch of transitions from the memory
32    to train the network. This helps in breaking the correlation between the inputs. I
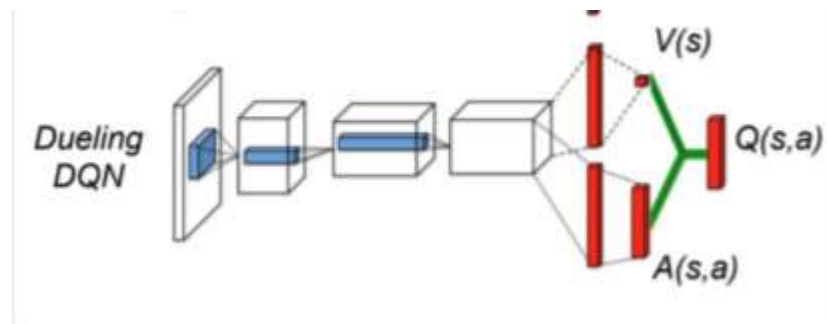33    used the python Deque collection to implement this.

34
35 ***Experience Replay***
36 • Policy changes rapidly with slight changes to Q-values, due to this, the
37 distribution of data can swing from one extreme to another. To overcome this,
38 we use two kinds of networks, the original network (w) and the target network
39 (w⁻). We use the target model to predict the Q-values. After C updates, we
40 synchronize w⁻ to w. By doing this, we fix the Q-value targets temporarily to
41 avoid a moving target.
42 • Scale of rewards and Q-values is unknown. This we have large rewards, the Q-
43 values increase a lot due to this to prevent this, we clip the rewards between -1
44 and 1. To do this, I used np.sign function.
45
46 **1.2   Dueling DQN:**
47 Q- Values tells us how good it is to take an action a in a state s. We can decompose
48 Q(s,a) as A(s,a) +V(s). Where A(s,a) is called Advantage function. If A>0 our
49 gradient is pushed in that direction. Else gradient is pushed in the opposite direction.
50 In dueling DQN, we separate the estimator using two new streams, one estimates the
51 state value V(S) and one estimates the advantage for each action A(s,a).



52

53 ***Dueling DQN***

54 While combining the two streams, we can directly add the two outputs:

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + A(s, a; \theta, \alpha)$$

55

56 But from this if we are given a Q, we cannot recover V and A uniquely which in turn
57 leads to poor practical performance.

58 Solutions:

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

59

60 This forces the advantage function to have zero advantage at chosen action.

61 We can also replace the max with the average:

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

This increases the stability of the optimizations. I implemented the average of the Advantage function in this code.

## 1.3 Reinforce

There are two broad ways of solving a reinforcement learning problem – They are calculating values functions like state values or Q-values by either approximating them or maintaining a tabular method. In this method, we approximate the state or action value to get the action and then use methods like epsilon decay to improve the policy. The second method is to directly approximate the policy directly.

$$\pi_\theta(s, a) = P[a|s, \theta]$$

*Policy approximation*

These are called policy gradient methods. They target at modelling and optimizing the policy directly. We perform gradient ascent here to maximize performance.

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$$

*Gradient Ascent*

In policy gradient we get the probability of each action and then choose action randomly. The outputs are calculated using the sigmoid activation function.

$$\pi(a|s, \theta) = \frac{e^{h(s,a,\theta)}}{\sum_b e^{h(s,b,\theta)}}$$

*Sigmoid activation to get probability of actions*

The finally gradient of J is written as

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi[Q^\pi(s, a)\nabla_\theta \ln \pi_\theta(a|s)]$$

*Gradient of J*

We use the log likelihood function to calculate the loss. Which is multiplied with the Return that we calculate this return can either be the return of the entire episode (MC policy gradient) or the Q-value calculated by the critic (1 Step TD Actor-Critic).

Reinforce is also called as Monte-Carlo Policy Gradient. In this, we calculate the total discounted return of the entire episode to calculate the gradient.

The update function looks like

$$\Delta\theta_t = \alpha G_t \nabla_\theta \log \pi_\theta(s_t, a_t)$$

*Reinforce update function*

99 Keras does not have a built-in loss function to handle log-likelihood loss. So, we
100 define a custom loss function that calculates the log likelihood of the predicted
101 values and the true values and multiplies them with the advantages and sums them
102 up.

103 The choose action function takes the probabilities from the model and chooses an
104 action based on the probabilities.

105 The learn function calculates the discounted sum of the rewards for the entire
106 episode. We use a baseline here that calculates the mean of the all rewards and
107 standard deviation of all the rewards and calculates sum-mean/standard deviation.
108 And train the model based on this baseline and states of the entire episode.

109

## 110 1.4 Actor-Critic

111

112 Reinforce has high variance to overcome this, we use a critic to estimate the action
113 value function along with the policy approximation. We can think of this as the actor
114 is the agent who does whatever action he wants and critic is tells the actor if that
115 action was good or not.

116

117 The actor updates theta and the critic updates w. The Q-values that the critic predicts
118 is sent to the actor using which the actor knows if this was a good action or not. The
119 actor uses a log likelihood loss similar to policy gradient and the actor uses mean
120 squared error as loss function.

121

122 The actor model takes state and delta(the TD error calculated using critic) and the
123 critic takes just the state as input and outputs the Q value for the state action pair.

124 The choose action function is similar to reinforce it gets the probabilities from the
125 actor and chooses a random action based on the probabilities.

126

127 The learn function calculates the critic values for the preset state and next state and
128 calculates the target and the td error. The actor is trained using this td error and the
129 critic is trained using the state and target.

## 130 1.5 PPO

131 In reinforce and other actor critic algorithm we still face a lot of variance in the
132 policy and drastic changes in the policy. The agent needs to take importance
133 sampling to reduce this variance. One way to do this is by clipping the ration of the
134 new policy and old policy by a certain factor. TRPO uses KL divergence to measure
135 the difference between the two new policies.

136 In PPO with clipped objectives, we maintain two networks, one with the current
137 policy that we want to refine and second that we use to collect samples.

$$\pi_\theta(a_t|s_t) \qquad \pi_{\theta_k}(a_t|s_t)$$

138

139 With the idea of importance sampling, we evaluate a policy samples collected from
140 older policy.

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t\left[\frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}\hat{A}_t\right]$$

141

142    We synchronize the second network to the first network using soft update. With
143    clipped objective, we calculate the ratio between the new policy and old policy. Our
144    objective function becomes:

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathop{\mathrm{E}}_{\tau \sim \pi_k} \left[ \sum_{t=0}^{T} \left[ \min(r_t(\theta)\hat{A}_t^{\pi_k}, \text{clip}\left(r_t(\theta), 1 - \epsilon, 1 + \epsilon\right)\hat{A}_t^{\pi_k}) \right] \right]$$

145

146    We are clipping the ration of the old and new policy between 1-€ and 1+€ generally
147    € is set to 0.2 in the PPO paper. This is very simple to implement.

148

149    **1.6    Environments:**

150

151    In this project, I have trained the agent in one environments for comparing the policy
152    rewards:

153

154    •    **LunarLander-v2:**

155

156    Landing pad is always at coordinates (0, 0). Coordinates are the first two numbers in state
157    vector. Reward for moving from the top of the screen to landing pad and zero speed is
158    about 100...140 points. If lander moves away from landing pad it loses reward back.
159    Episode finishes if the lander crashes or comes to rest, receiving additional -100 or +100
160    points. Each leg ground contact is +10. Firing main engine is -0.3 points each frame.
161    Solved is 200 points. Landing outside landing pad is possible. Fuel is infinite, so an agent
162    can learn to fly and then land on its first attempt. Four discrete actions available: do
163    nothing, fire left orientation engine, fire main engine, and fire right orientation engine.

164

165    *Possible Actions:*

166

167    Actions here is four floats [main engine, left-right engines, fire main engine, and fire
168    right orientation engine]

169

170    •    **Acrobot-v1:**

171

172    Acrobot is a 2-link pendulum with only the second joint actuated. Initially, both links
173    point downwards. The goal is to swing the end-effector at a height at least the length of
174    one link above the base. Both links can swing freely and can pass by each other, i.e., they
175    don't collide when they have the same angle.

176

177    *State:*

178

179    The state consists of the sin () and cos () of the two rotational joint angles and the joint
180    angular velocities:
181    [cos (theta1) sin (theta1) cos (theta2) sin (theta2) thetaDot1 thetaDot2]. For the first link,
182    an angle of 0 corresponds to the link pointing downwards. The angle of the second link is
183    relative to the angle of the first link. An angle of 0 corresponds to having the same angle
184    between the two links. A state of [1, 0, 1, 0…] means that both links point downwards.

185

*Actions:*

187

188     The action is either applying +1, 0 or -1 torque on the joint between the two pendulum
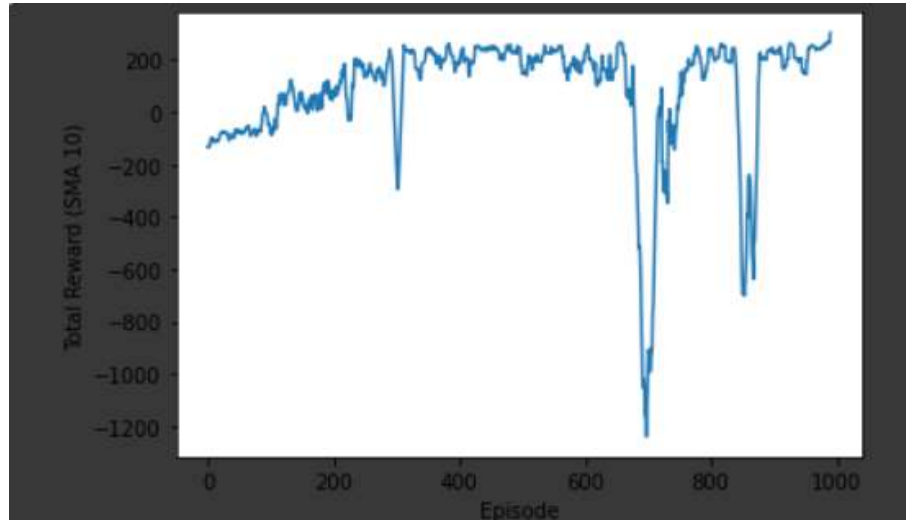189     links.

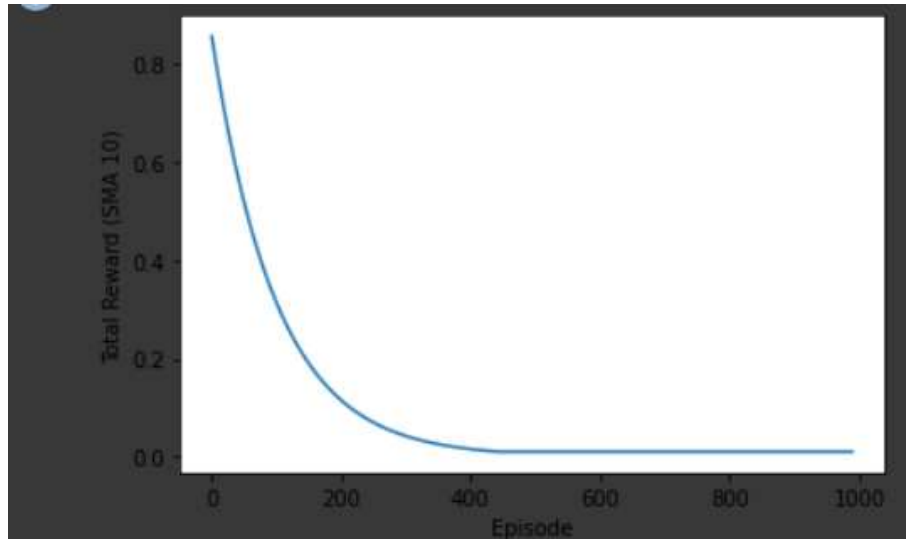190     **2      Results:**

191

192     **2.1     DQN:**

193

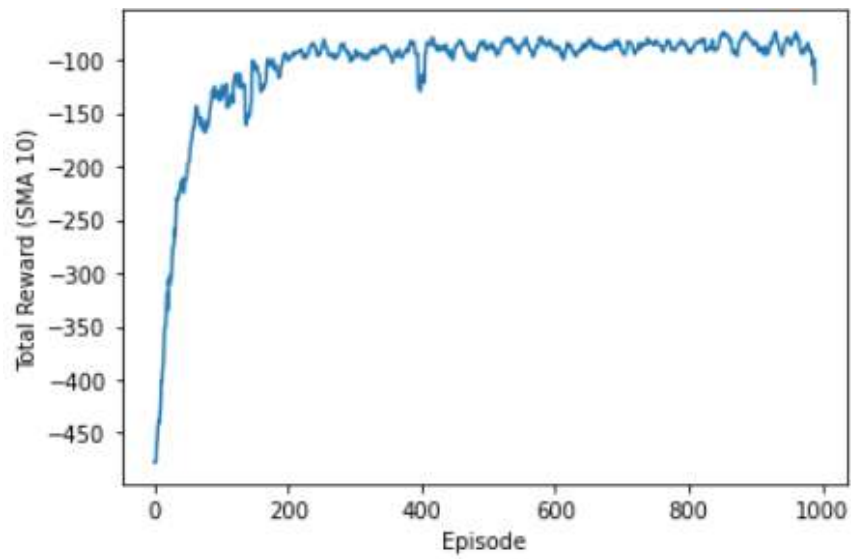194          • **LunarLander-v2**



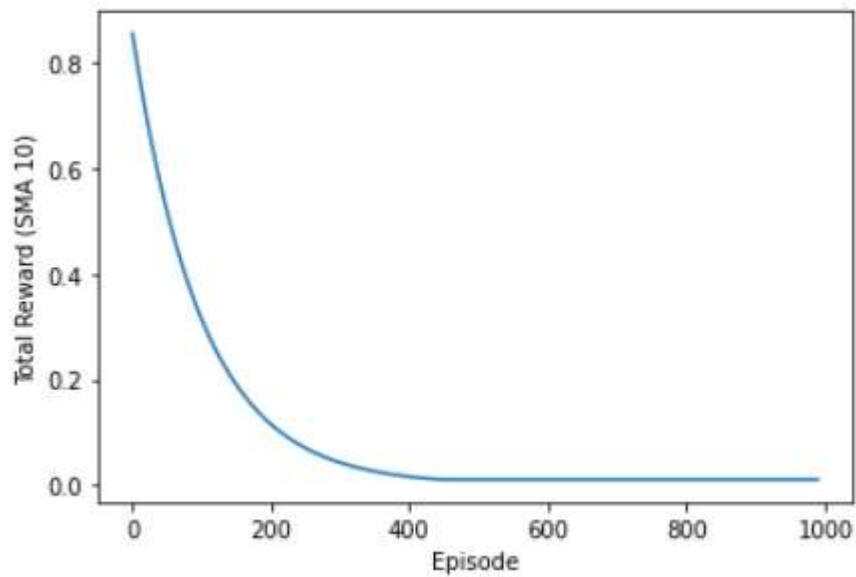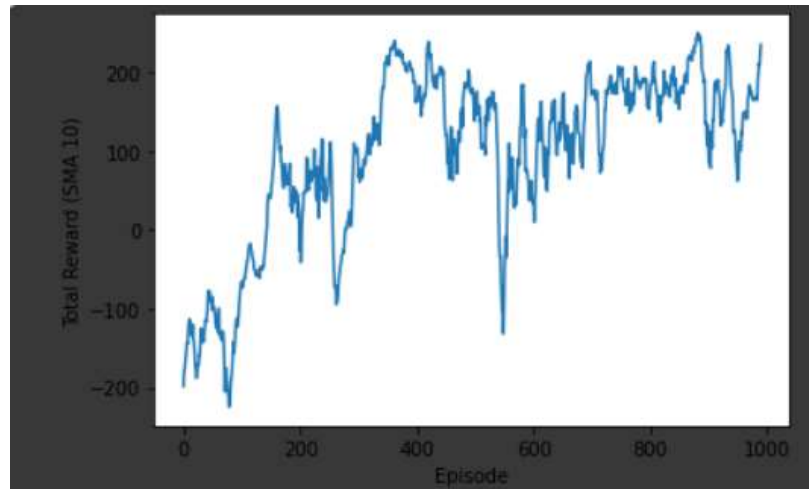195



196
197

198          • **Acrobot-v1**
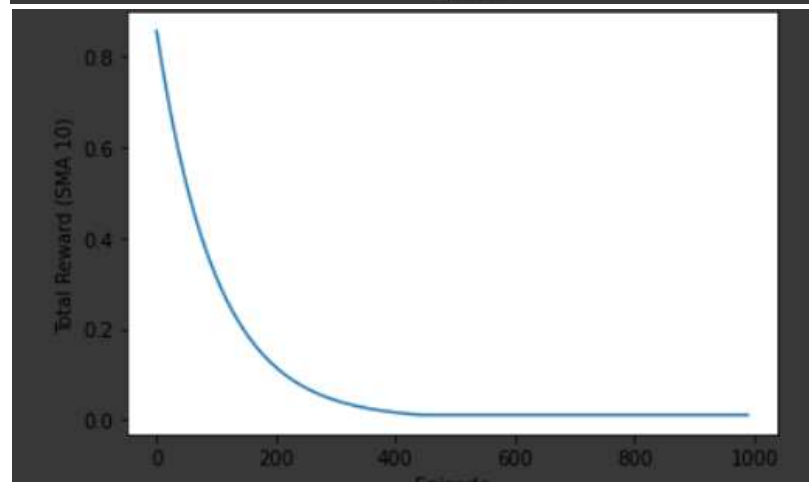
199

200
201
202
203    **2.2      Dueling DQN:**
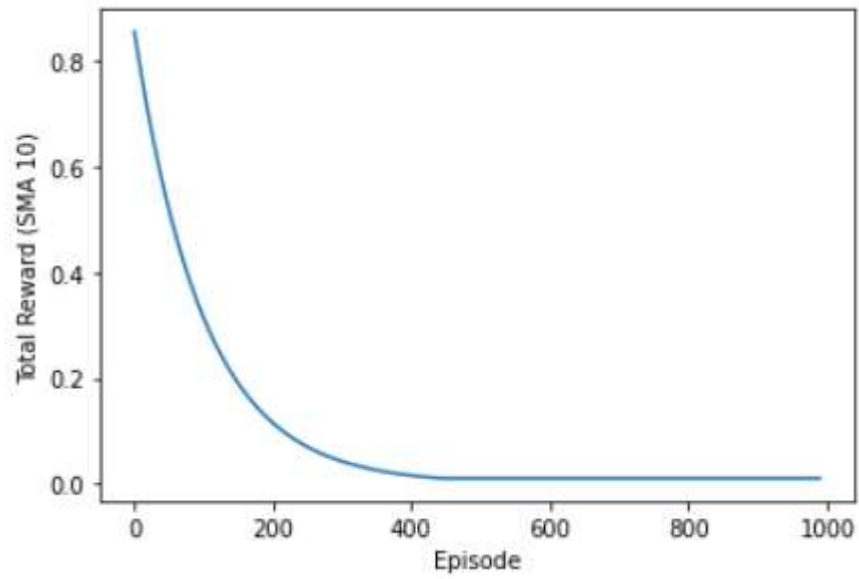204
205            • **LunarLander-v2**

206     •

207     •

208

209     •   **Acrobot-v1**



210

211
212
213 **2.3**    **Reinforce:**
214
215        • **LunarLander-v2**



216
217        • **Acrobot-v1**
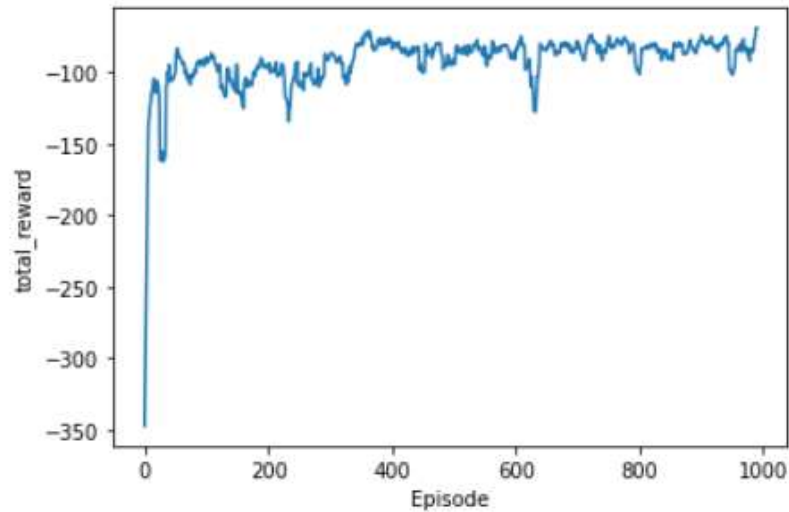
218

219

220

221 **2.4** **Actor-Critic:**

222

223  • **LunarLander-v2**



224

225  • **Acrobot-v1**

226

227 **2.5** **PPO:**



228

229

230 **3** **Conclusion:**

231 From the above graphs we can infer the following thing:

- DQN trains very fast somewhere around 200 episodes but it isn't as stable as dueling DQN. DQN overestimates a few actions.
- DQN oscillates to negative reward suddenly and then comes back the max reward whereas dueling DQN changes a lot but only in the positive reward region.
- Policy gradients algorithms like REINFORCE and 1 step AC need a lot of samples to converge.
- REINFORCE has a lot of variance as seen from the graph and AC slightly decreases the variance.
- PPO is by far the most stable Actor-Critic algorithm but it needs lot of samples to converge.
- All of the policy gradient algorithms and AC critic algorithms are highly

sensitive to hyper parameter tuning. I am sure we can get better results for some different set of hyper parameters.

## 4    References

1) https://gym.openai.com/envs/LunarLander-v2/

2) Lecture slides – CSE510 Introduction to Reinforcement Learning

3)https://medium.com/@jonathan_hui/rl-proximal-policy-optimization-ppo-explained-77f014ec3f12

4) Human-level control through deep reinforcement learning (Mnih et. al)

5) Dueling Network Architectures for Deep Reinforcement Learning (Wang et.al)

6)https://towardsdatascience.com/understanding-actor-critic-methods-931b97b6df3f?gi=699deab97419

7) Proximal Policy Optimization Algorithms https://arxiv.org/pdf/1707.06347.pdf