

Controlling 3D gaming agents in an adversarial setting with Deep Reinforcement Learning

Mehmood Munir

p176075@nu.edu.pk

Bashir Ahmed

p176079@nu.edu.pk

M. Hanzaila

p180453@nu.edu.pk

Supervisor

Dr. Muhammad Nauman

January 6, 2021

Table of contents

1. Q-Learning
2. Calculation
3. Bridge Design Pattern
4. Coding Convention
5. Airstriker
6. CartPole
7. Slimevolleygym
8. Work Breakdown

Q Learning

Q Learning

$$\underbrace{Q(s, a)}_{\text{New Q-Value}} = Q(s, a) + \underbrace{\alpha}_{\text{Learning Rate}} \left[\underbrace{R(s, a, s')}_{\text{Reward}} + \underbrace{\gamma}_{\text{Discount rate}} \underbrace{\max_{a'} Q'(s', a')}_{\text{Maximum predicted reward, given new state and all possible actions}} - Q(s, a) \right]$$

$$\text{Sample} = R(s, a, s') + \gamma(\max_{a'} Q(s', a'))$$

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha(\text{Sample})$$

Grid World

Grid World

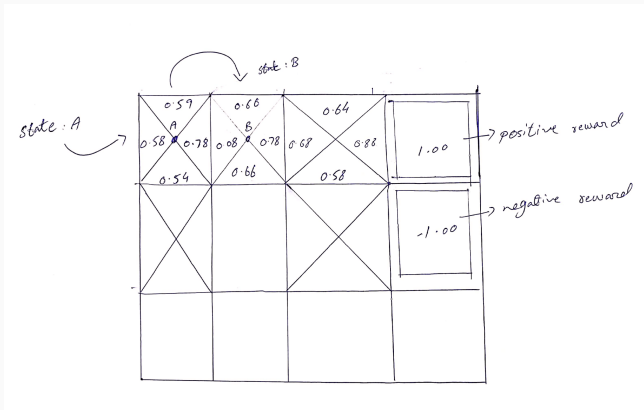


Figure 1: Grid World

Calculation

Calculation

$$R(s,a,s') = 1$$

s = state(A),

s' = state(B) next state,

a = move right(Random Action)

Q-Values with different actions 0.66, 0.40, 0.78, 0.66

$$\max Q(s',a') = \max (0.66, 0.40, 0.78, 0.66) = 0.78$$

$$\text{Sample} = R(s,a,s') + \gamma(\max Q(s', a'))$$

$$\text{Sample} = 1 + 1(0.78) = 1.78$$

$$Q(s,a) = (1-\alpha)Q(s, a) + \alpha(\text{Sample})...(A)$$

Let $\alpha = 1$

putting values in equation (A)

$$Q(s,a) = (1 - 1)(0.78) + 1(1.78)$$

$$Q(s,a) = 1.78$$

Bridge Design Pattern

Bridge Design Pattern

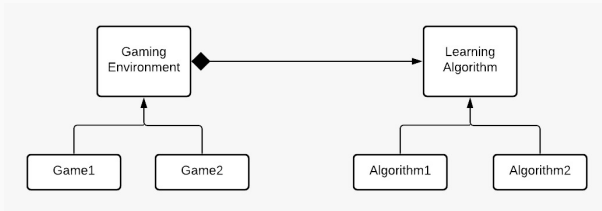


Figure 2: Factory Design Pattern

Coding Convention

- Tensorflow v2 vs tensorflow v1

Coding Convention

- Tensorflow v2 vs tensorflow v1
- Snake case naming convention

Coding Convention

- Tensorflow v2 vs tensorflow v1
- Snake case naming convention

Airstriker

- State Space: 215040
- Action Space : 12

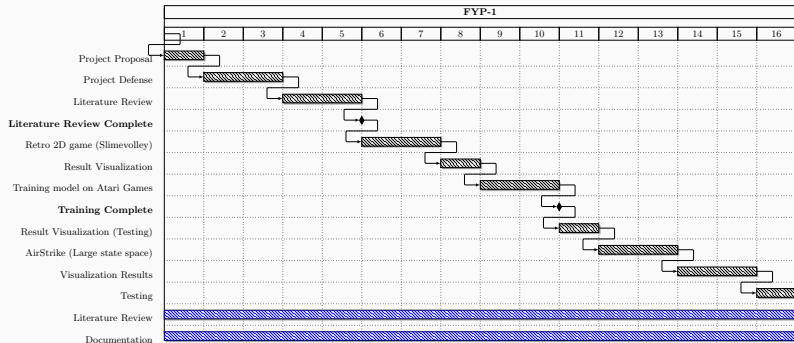
CartPole

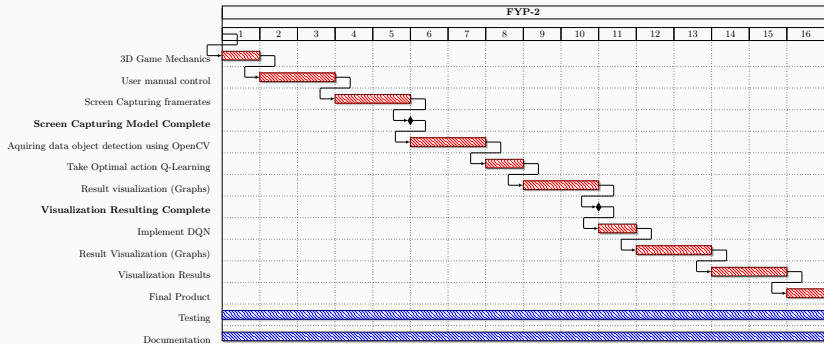
- State Space: 16
- Action Space : 2

Slimevolleygym

- State Space: 12
- Action Space : 3

Work Breakdown





Questions?