# CS529– Applied Artificial Intelligence
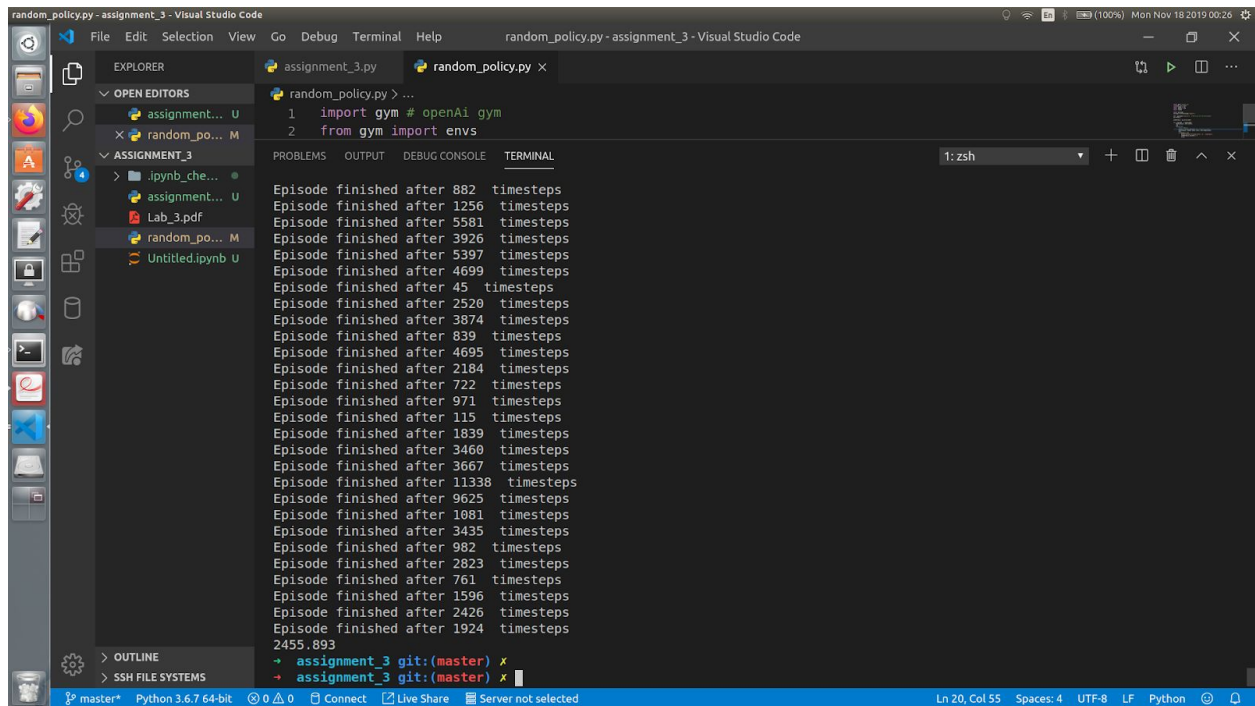# Lab Assignment - 3

**Name** : Shah Jainam Mukeshbhai
**Entry Number** : 2017csb1107

## Question 1 :

Chose an action randomly and when the reward was **20**, then reseted the environment. The average number of timesteps for 1000 episodes were **2455.893**(Around 2500).



## Question 2 :

| Discount Factor | Policy Iteration | Value Iteration |
|---|---|---|
| 0.99 | 12 * 1100 | 725 |
| 0.95 | 12 * 226 | 146 |
| 0.9 | 12 * 110 | 74 |

| 0.8 | 12 * 53 | 38 |
|-----|---------|-----|
| 0.4 | 12 * 14 | 14 |

- For discount factor 0.4, I didn't get the same policy
- Otherwise, Optimal policy is same for Value Iteration and Policy Iteration
- Compared to each other, value-iteration is computationally efficient even though it takes more number of iterations to converge, each iteration is less computationally expensive than policy-iteration.
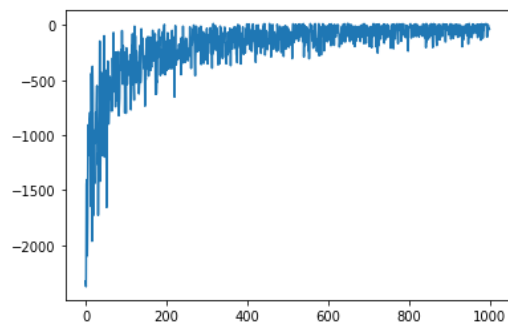
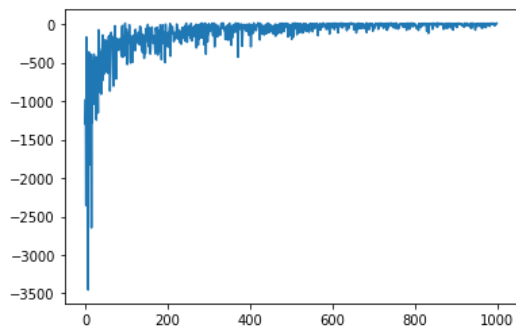Distribution of Number of steps taken for 1000 episodes is as follows:



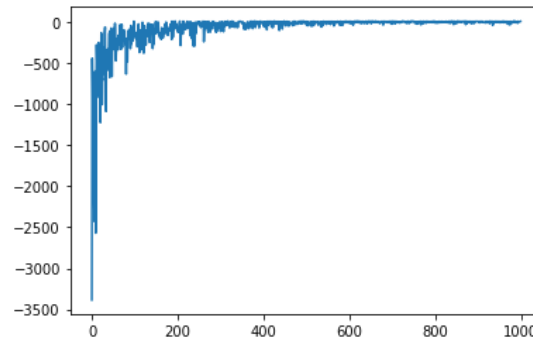Average number of steps is 13.

## Question 3 :

Alpha = 0.05, Convergence around greater than 1000 eps

Alpha = 0.1, Convergence around 600 eps                Alpha = 0.2, Convergence around 500 eps
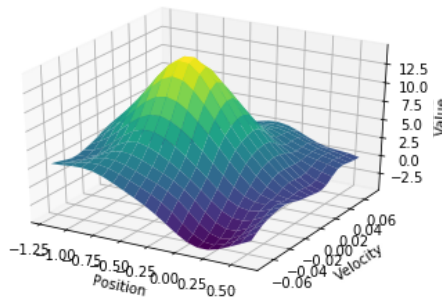
- On Y-axis, I kept total reward of the episode.
- On convergence, the change is total reward becomes constant. It becomes parallel to x-axis.
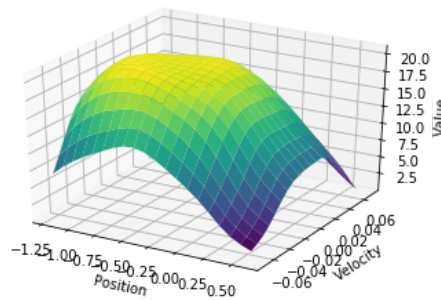
# Question 4 :

The results were matching the graphs of book. Below are the plots of Mountain Car problem.

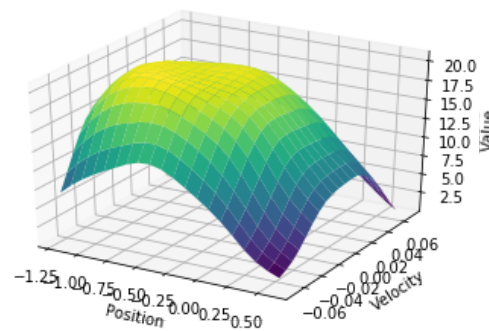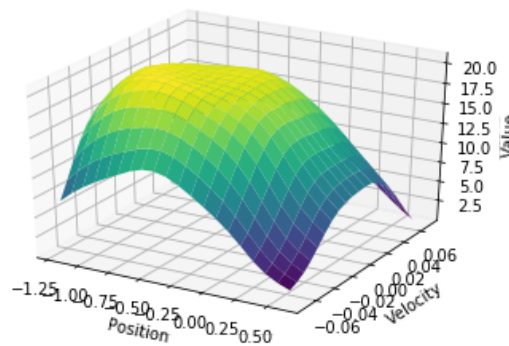Episode 0                                                Episode 100
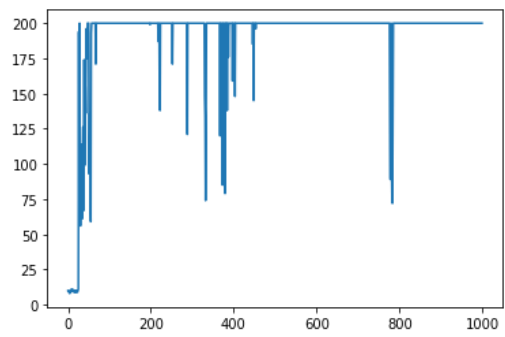


Episode 200                                              Episode 300

For Cart Pole, the plot of Total reward with num of episodes is :

References:
For Mountain Car:
https://github.com/SamKirkiles/mountain-car-SARSA-AC/blob/master/mountain_car.py
For Cart Pole: https://github.com/ceteke/RL/blob/master/Approximation/Linear%20Sarsa.ipynb