# Assignment -2

## OpenAI GYM:



## QLearning:

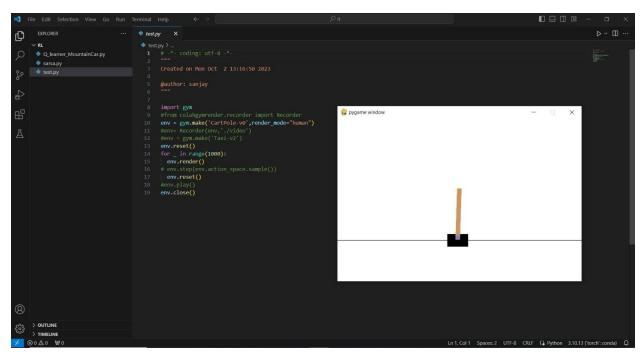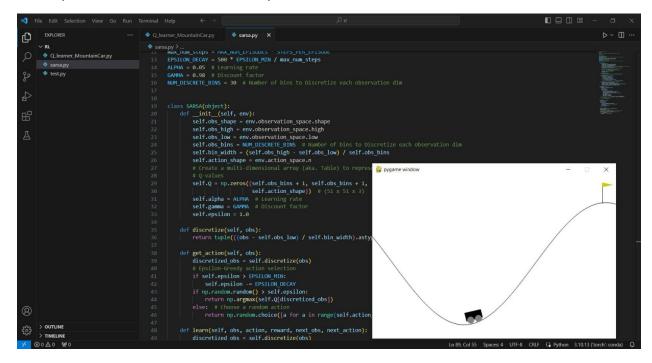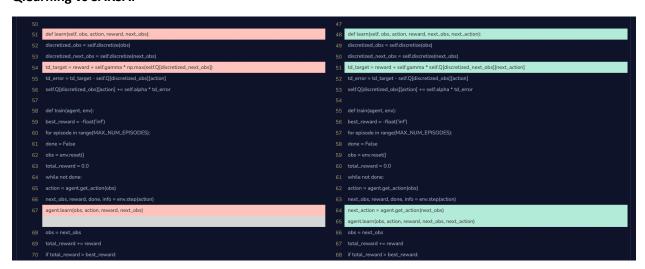**SARSA (State Action Reward State Action):**



**Qlearning vs SARSA:**



GIT: https://github.com/sanjaybhargavm/mountain_car.git