

Sumeet Shanbhag
Student ID : 641020714

Reinforcement Learning Monte Carlo **Programming Assignment**

Output:

Below is the output for optimal action value functions and optimal policies using the exploring start and first visit methods after running for 500 episodes.

```
[Running] python -u "c:\Users\Sumeet\OneDrive - Worcester Polytechnic Institute (wpi.edu)\WPI Sem 3\Reinforcement Learning\Week 4\HW.py"
Exploring Starts

optimal action value function:
{(0, 'T'): 1.0, (1, 'L'): 1.6676888447942863, (1, 'R'): 3.7334822137842147, (2, 'L'): 3.6580228324669079, (2, 'R'): 4.079527295188684, (3, 'L'): 3.942815791136172, (3, 'R'): 4.39206375321582, (4, 'L'): 4.217593556601563, (4, 'R'): 4.67380313229167, (5, 'T'): 5.0}

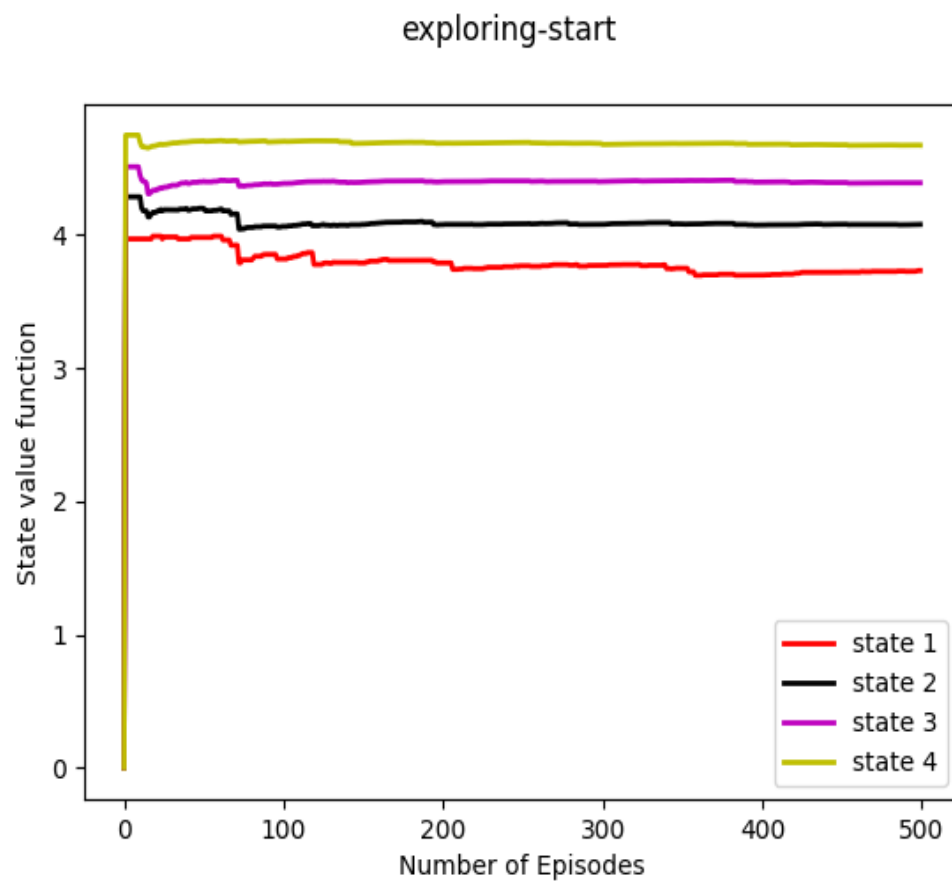
optimal policy:
[(0, 'T'), (1, 'R'), (2, 'R'), (3, 'R'), (4, 'R'), (5, 'T')]

First Visit

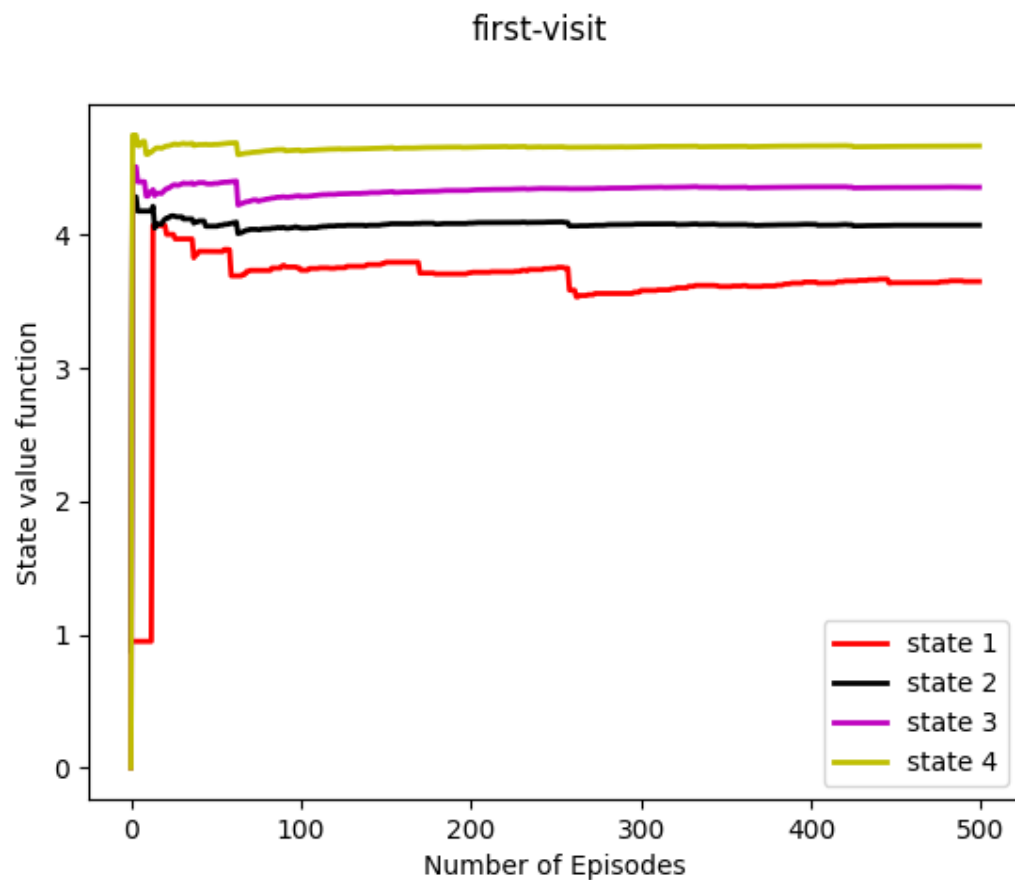
optimal action value function:
{(0, 'T'): 1.0, (1, 'L'): 1.3830245805884593, (1, 'R'): 3.651488112084768, (2, 'L'): 3.496311737894315, (2, 'R'): 4.07381866904414, (3, 'L'): 3.9961351813216575, (3, 'R'): 4.358595784792832, (4, 'L'): 4.248734446188556, (4, 'R'): 4.668751498171741, (5, 'T'): 5.0}

optimal policy:
[(0, 'T'), (1, 'R'), (2, 'R'), (3, 'R'), (4, 'R'), (5, 'T')]
```

The figure below shows the estimation of the state value function vs the number of episodes it takes for Monte Carlo method of exploring starts



The figure below shows the estimation of the state value function vs the number of episodes it takes for Monte Carlo method of On-policy first visit



Convergence:

Exploring starts method ensures every state-action pair is sampled at least once. This guarantees complete exploration of the environment, allowing for faster and more accurate convergence of state value estimations, especially when compared to methods that might miss certain state-action pairs.

The first-visit Monte Carlo method estimates state values by averaging returns from the first occurrence of each state in episodes. While it provides a consistent estimate as the number of episodes increases, its convergence can be slower, especially for states that are infrequently visited.