

EX5

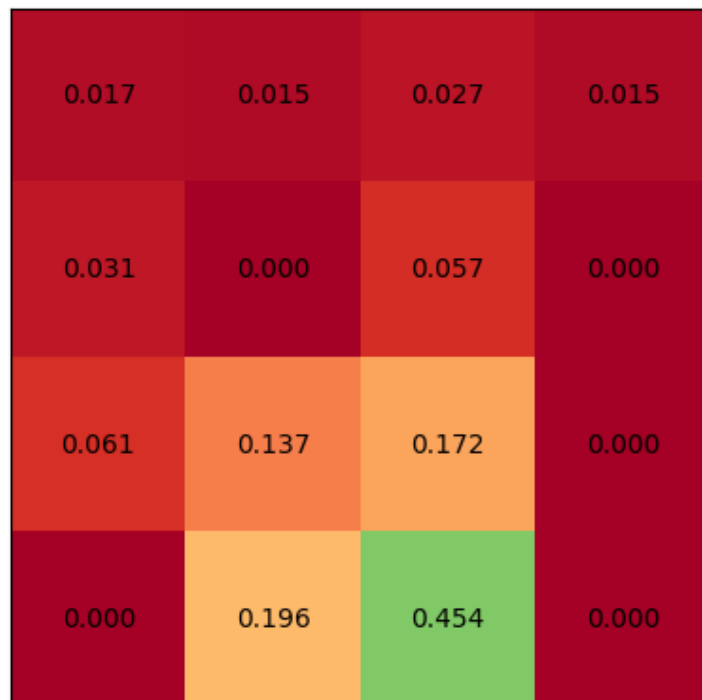
Random Walk

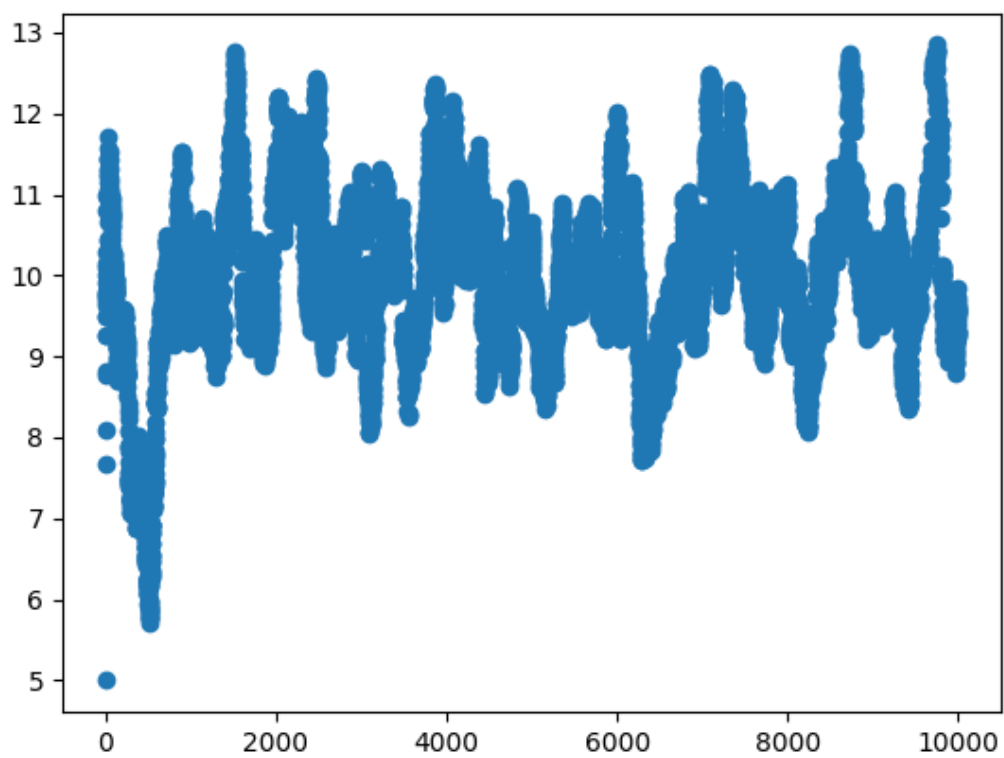
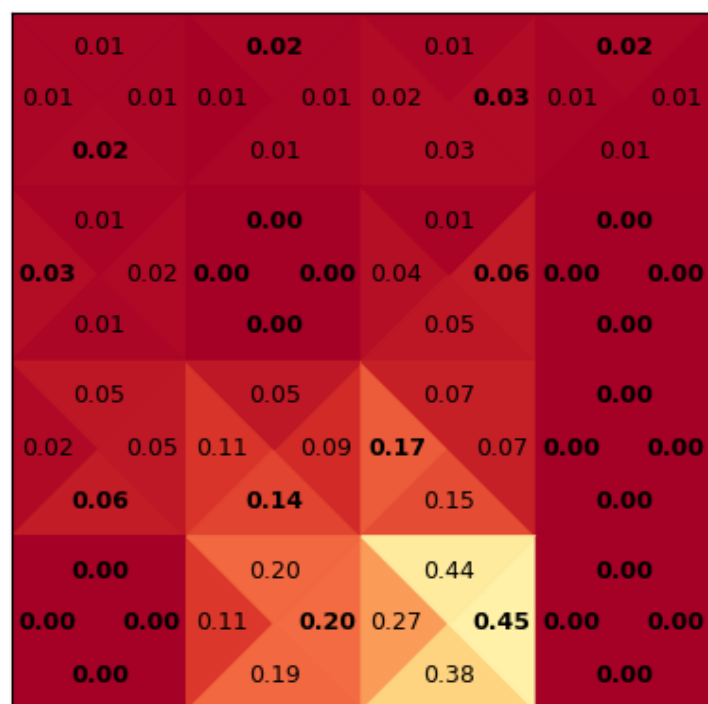
The random walk terminated in the state to the left, what it did before cannot be determined, except we know it did not jump from state E to the right.

Only the estimate for state A was changed, because with a discount factor of 1, the TD-error was zero for all other states.

$$\begin{aligned} V(A) &= V(A) + 0,1 (0 + 1 \cdot 0 - 0,5) \\ &= 0,5 - 0,05 \\ &= 0,45 \end{aligned}$$

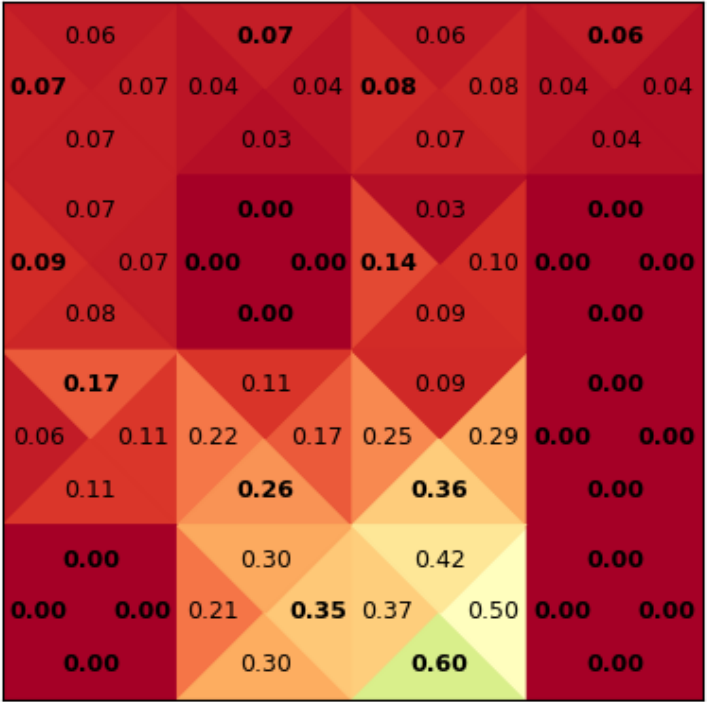
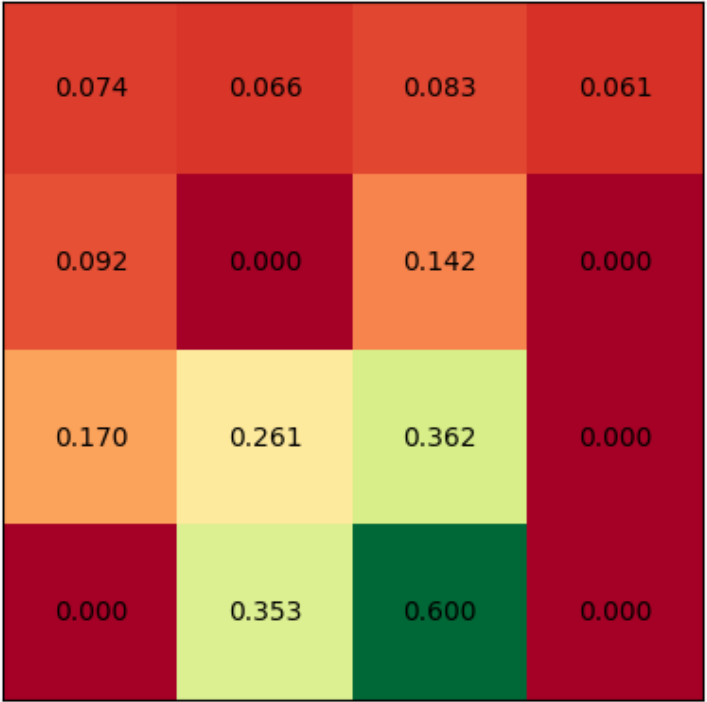
Sarsa

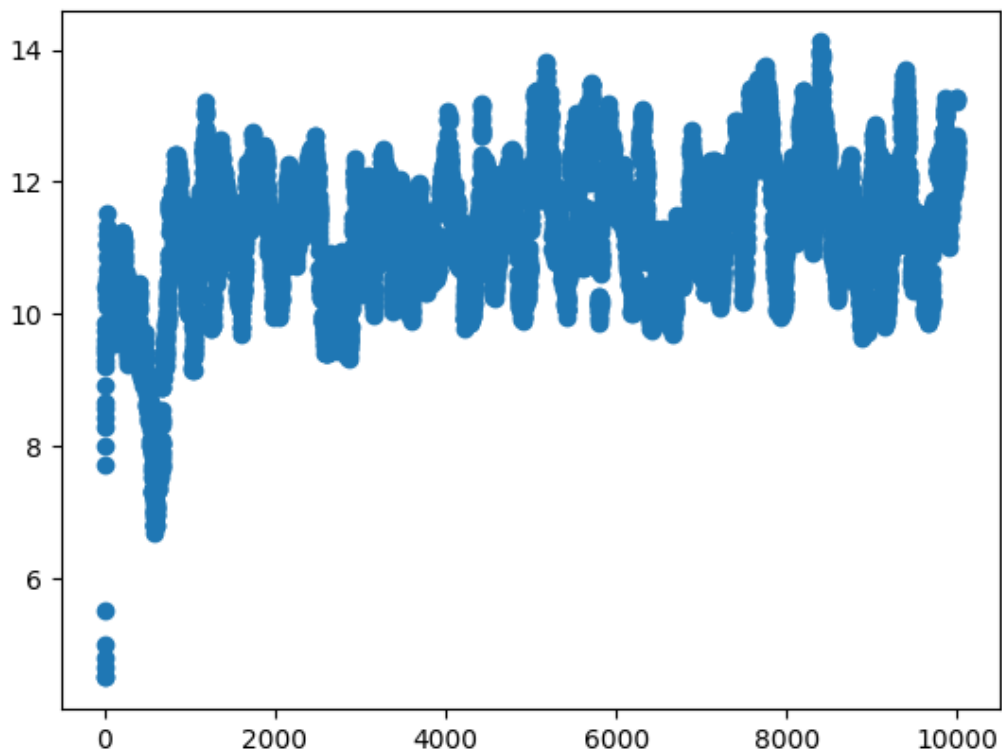




The average episode length does not change during training. The total average episode length is 9,997

Q-Learning





The average episode length for Q-learning is slightly longer compared to SARSA. The total average is 11,2754

The console output for the small stochastic environment is:

current environment:

SFFF

FHFH

FFFH

HFFG

Running sarsa...

↓ ↑ → ↑

← H → H

↓ ↓ ← H

H → → G

episode_length: 9.997

Running qlearning

← ↑ ← ↑

← H ← H

↑ ↓ ↓ H

H → ↓ G

episode_length: 11.2754

And for the deterministic environment:

current environment:

SFFF

FHFH

FFFH

HFFG

Running sarsa...

↓ → ↓ ←

↓ H ↓ H

→ → ↓ H

H → → G

episode_length: 7.772

Running qlearning

↓ → ↓ ←

↓ H ↓ H

→ ↓ ↓ H

H → → G

episode_length: 7.3249

When changing to the deterministic environment, both algorithms still not create the same policy, but the average episode length drops significantly.

The larger environment exhibits the same behaviour as the small one.