

Question 8 Exercise 6.12

If action selection was greedy, they will act similar.

The action selection at every episode step would be the same.

The SARSA update equation essentially becomes

Choose A' such that ~~Q~~ we get $\max_a(Q(s, A'))$

$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s, A') - Q(s, A)]$$

$$= Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_a Q(s, a) - Q(s, A)]$$

which is the same as the Q learning update.

Updates and action selection would be the same for both algorithms provided a greedy selection replaces all ϵ -greedy policies involved.