Felix Bühler - 2973140
Jan Leusmann - 2893121
Jamie Ullerich - 3141241

# Exercise 5

## Task 1 - Random Walk

- The first episode terminates in the left state.

- The states are not yet visited or the ones visited are updated with TD error $= 0$
  ($\rightarrow$ the values are not changing at all.)

- $V_{t+1}(A) = V_t(A) + \alpha * (R_{t+1} + \gamma V(terminal) - V_t(A))$ in our case: $0.5 + 0.1 * (0 + 0 - 0.5) = 0.45$

## Task 2 - Sarsa and Q-learning on the FrozenLake

**a)**

see Figure 1 & Figure 2 & Figure 3.

$$
\begin{array}{cccc}
\leftarrow & \uparrow & \downarrow & \uparrow \\
\leftarrow & \leftarrow & \leftarrow & \leftarrow \\
\uparrow & \downarrow & \leftarrow & \leftarrow \\
\leftarrow & \rightarrow & \uparrow & \leftarrow
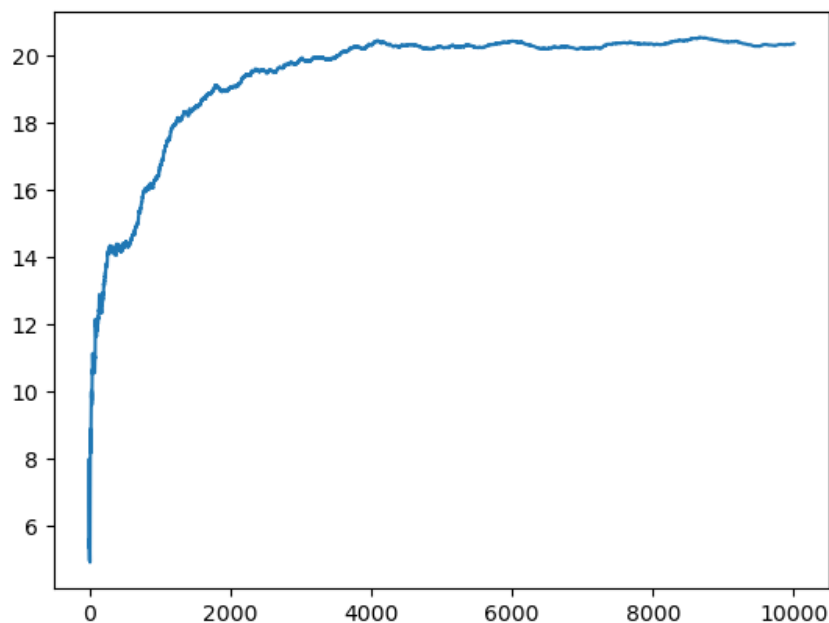\end{array}
$$

Table 1: Sarsa policy



Figure 1: Sarsa training length

Figure 2: Sarsa V



Figure 3: Sarsa Q

**b)**

Sarsa follows a more saver policy (more exploration) then q-learning (more exploitation). Same as in the cliff example in the lecture Q-learning takes the optimal path which could end up in holes.

see Figure 4 & Figure 5 & Figure 6.

| | | | |
|---|---|---|---|
| ↓ | ↑ | ← | ↑ |
| ← | ← | ← | ← |
| ↑ | ↓ | ← | ← |
| ← | → | → | ← |

Table 2: Q-Learning policy

Figure 4: Q-Learning training length



Figure 5: Q-Learning V

**c)**

An optimal policy gets calculated. This policy would not be optimal in a slippery environment.
    For Sarsa: see Figure 7 & Figure 8 & Figure 9.

$$
\begin{array}{cccc}
\downarrow & \rightarrow & \downarrow & \leftarrow \\
\downarrow & \leftarrow & \downarrow & \leftarrow \\
\rightarrow & \downarrow & \downarrow & \leftarrow \\
\leftarrow & \rightarrow & \rightarrow & \leftarrow
\end{array}
$$

Table 3: det Sarsa policy

For Q-Learning: see Figure 10 & Figure 11 & Figure 12.

Figure 6: Q-Learning Q



Figure 7: det Sarsa training length

| | | | |
|---|---|---|---|
| ↓ | → | ↓ | ← |
| ↓ | ← | ↓ | ← |
| → | → | ↓ | ← |
| ← | → | → | ← |

Table 4: det Q-Learning policy

## d)

For Sarsa: see Figure 7 & Figure 8 & Figure 9.

For Q-Learning: see Figure 16 & Figure 17 & Figure 18.
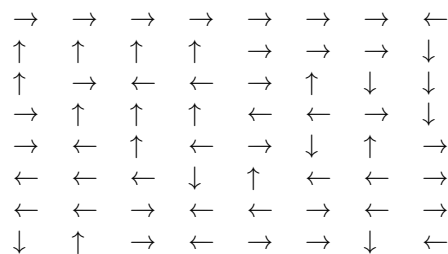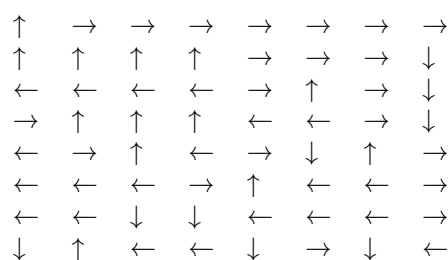
Figure 8: det Sarsa V



Figure 9: det Sarsa Q

| → | → | → | → | → | → | → | ← |
|---|---|---|---|---|---|---|---|
| ↑ | ↑ | ↑ | ↑ | → | → | → | ↓ |
| ↑ | → | ← | ← | → | ↑ | ↓ | ↓ |
| → | ↑ | ↑ | ↑ | ← | ← | → | ↓ |
| → | ← | ↑ | ← | → | ↓ | ↑ | → |
| ← | ← | ← | ↓ | ↑ | ← | ← | → |
| ← | ← | → | ← | ← | → | ← | → |
| ↓ | ↑ | → | ← | → | → | ↓ | ← |

Table 5: 8x8 Sarsa policy

| ↑ | → | → | → | → | → | → | → |
|---|---|---|---|---|---|---|---|
| ↑ | ↑ | ↑ | ↑ | → | → | → | ↓ |
| ← | ← | ← | ← | → | ↑ | → | ↓ |
| → | ↑ | ↑ | ↑ | ← | ← | → | ↓ |
| ← | → | ↑ | ← | → | ↓ | ↑ | → |
| ← | ← | ← | → | ↑ | ← | ← | → |
| ← | ← | ↓ | ↓ | ← | ← | ← | → |
| ↓ | ↑ | ← | ← | ↓ | → | ↓ | ← |

Table 6: 8x8 Q-Learning policy

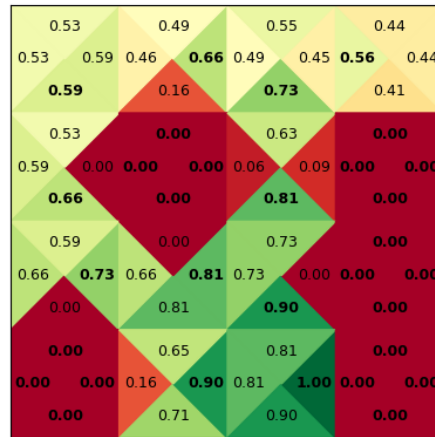Figure 10: det Q-Learning training length



Figure 11: det Q-Learning V

Figure 12: det Q-Learning Q
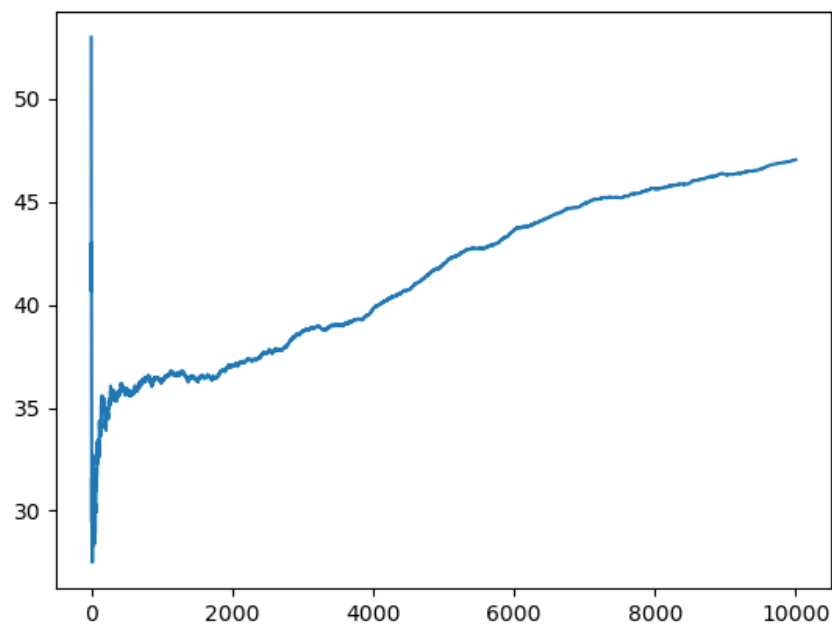


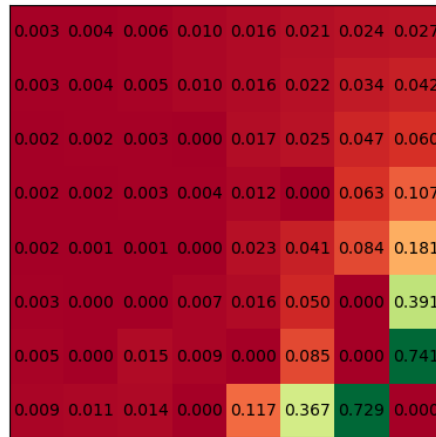Figure 13: 8x8 Sarsa training length
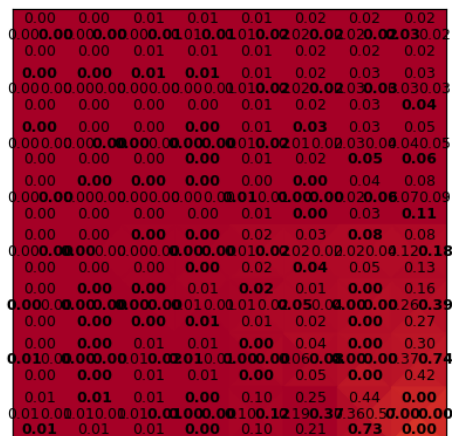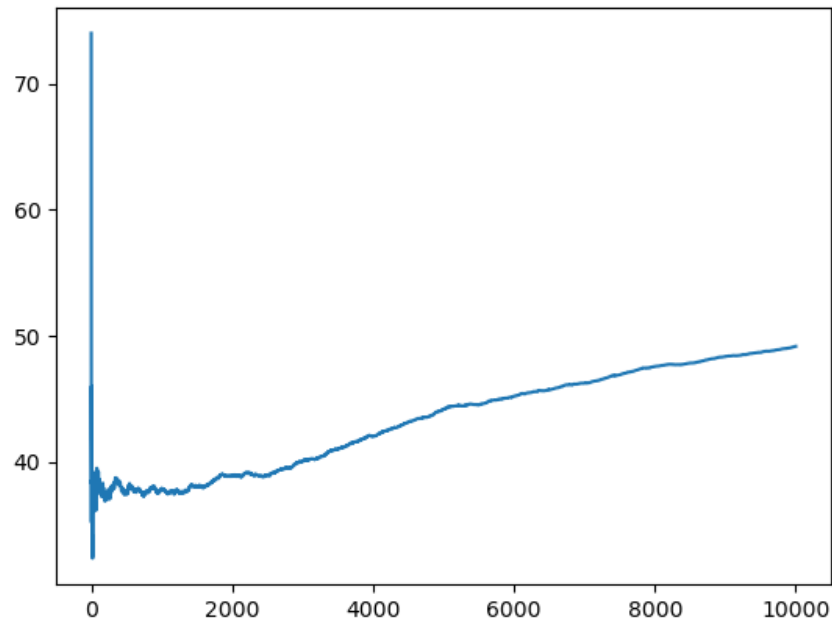
Figure 14: 8x8 Sarsa V
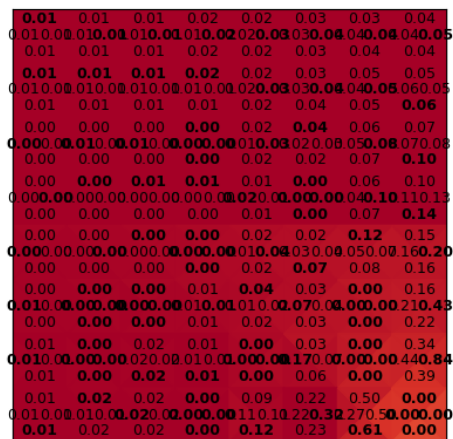


Figure 15: 8x8 Sarsa Q

Figure 16: 8x8 Q-Learning training length



Figure 17: 8x8 Q-Learning V

Figure 18: 8x8 Q-Learning Q