# Deep Reinforcement Learning S2

Dynamic Programming, Monte-Carlo Sarsa, Q-learning

Neta Bar-gil 305609836 & Yair Lahad 205493018

1. Dynamic Programming:
   For the Dynamic Programming optimal solution for the Frozen lake, we decided to implement Value iteration model using synchronous dynamic programming. We provide the MDP model that includes PSS' matrix for each action, and the given reward for each state. (Appendix 1) The optimal solution described in figures bellow:

   Optimal state values:

   Optimal policy:

   

2. Monte-Carlo control:
   In order to find the best approximation to the optimal solution found using the DP, we tried various parameters and compared the optimal state values from MC to the optimal state values received from the DP. Each parameter was tested for 500 episodes.

   Tested Parameters:

   **Epsilon = 0.9**

   Constant Alpha = [**0.01**,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]

   Epsilon decay functions = [**1/((n/10)+1),** 0.9*epsilon]

   The optimal solution (=best approximation) received by implementing the values appears in bold. After finding the optimal parameters we run the algorithm for 5000 episodes.

   The optimal solution for MC control described in figures bellow:

Optimal state values:    Optimal Action state values:    Optimal policy:

gridworld

| | 0.652 | 0.737 | 0.943 | |
| 1 | 5 | 9 | 13 | 17 |
| | 0.28 | | 0.884 | 0.945 |
| 2 | 6 | 10 | 14 | 18 |
| -0.024 | 0.204 | 0.486 | 0.789 | 0.873 |
| 3 | 7 | 11 | 15 | 19 |
| -0.041 | 0.016 | | 0.154 | |
| 4 | 8 | 12 | 16 | 20 |

gridworld (Optimal Action state values)

| | 0.008 / -0.071 0.652 0.01 / 0.002 | 0.013 / 0.737 / -0.06 | 0.058 / 0.017 0.943 / 0.031 | |
| 1 | 5 | 9 | 13 | 17 |
| | 0.28 / -0.063 -0.12 / -0.054 | | -0.01 / 0.01 0.884 0.0 / -0.021 | 0.945 / 0.009 / -0.013 |
| 2 | 6 | 10 | 14 | 18 |
| -0.257 / -0.256 -0.024 -0.232 / -0.242 | 0.204 / -0.156 -0.077 / -0.166 | -0.089 / 0.486 -0.016 / -0.073 | 0.011 / 0.789 0.008 / -0.013 | 0.873 / 0.008 / 0.008 |
| 3 | 7 | 11 | 15 | 19 |
| -0.418 / -0.358 -0.041 -0.254 / -0.359 | 0.016 / -0.275 / -0.228 | | 0.154 / 0.01 -0.012 / -0.015 | |
| 4 | 8 | 12 | 16 | 20 |

gridworld (Optimal policy)

| | → | → | → | |
| 1 | 5 | 9 | 13 | 17 |
| | ↑ | | → | ↑ |
| 2 | 6 | 10 | 14 | 18 |
| → | ↑ | → | → | |
| 3 | 7 | 11 | 15 | 19 |
| → | ↑ | | ↑ | |
| 4 | 8 | 12 | 16 | 20 |

we can see that this approximation is almost converging to the optimal solution, and there for will suggest increasing number of episodes (CPU limitation)

3. SARSA:
In order to find the best approximation to the optimal solution found using the DP, we tried various parameters and compared the optimal state values from SARSA to the optimal state values received from the DP. Each parameter was tested for 500 episodes.

Tested Parameters:

**Epsilon = 0.9**

Constant Alpha = [0.01,0.1,0.2,**0.3**,0.4,0.5,0.6,0.7,0.8,0.9,1.0]

Epsilon decay functions = [**1/((n/10)+1)**, 0.9*epsilon]

The optimal solution (=best approximation) received by implementing the values appears in bold. After finding the optimal parameters we run the algorithm for 5000 episodes.

The optimal solution for SARSA described in figures bellow:

Optimal state values:    Optimal Action state values:    Optimal policy:

gridworld

| | -0.398 | -0.385 | -0.251 | |
| 1 | 5 | 9 | 13 | 17 |
| | -0.4 | | -0.385 | -0.251 |
| 2 | 6 | 10 | 14 | 18 |
| -0.4 | -0.4 | -0.399 | -0.398 | -0.385 |
| 3 | 7 | 11 | 15 | 19 |
| -0.4 | -0.4 | | -0.41 | |
| 4 | 8 | 12 | 16 | 20 |

gridworld (Optimal Action state values)

| | -0.466 / -0.472 -0.398 -0.405 / -0.407 | -0.385 / 0.405 -0.463 / -0.468 | -0.403 / -0.4 -0.25 / -0.399 | |
| 1 | 5 | 9 | 13 | 17 |
| | -0.4 / -0.475 -0.471 / -0.47 | | -0.462 / -0.459 -0.385 -0.401 / -0.456 | -0.251 / -0.4 / -0.407 |
| 2 | 6 | 10 | 14 | 18 |
| -0.469 / -0.405 -0.405 -0.41 / -0.4 | -0.411 / -0.4 -0.47 / -0.407 | -0.471 / -0.399 -0.41 / -0.471 | -0.398 / -0.411 -0.404 / -0.409 | -0.385 / -0.436 / -0.404 |
| 3 | 7 | 11 | 15 | 19 |
| -0.4 / -0.4 -0.403 -0.4 / -0.402 | -0.469 / -0.405 / -0.405 | | -0.41 / -0.47 -0.475 / -0.476 | |
| 4 | 8 | 12 | 16 | 20 |

gridworld (Optimal policy)

| | → | ↑ | → | |
| 1 | 5 | 9 | 13 | 17 |
| | ↑ | | → | ↑ |
| 2 | 6 | 10 | 14 | 18 |
| ↓ | → | → | ↑ | ↑ |
| 3 | 7 | 11 | 15 | 19 |
| ← | ← | | ↑ | |
| 4 | 8 | 12 | 16 | 20 |

we can see that this approximation is not converging to the optimal solution, and there for will suggest increasing number of episodes (CPU limitation)

4. Q learning:
   In order to find the best approximation to the optimal solution found using the DP, we tried various parameters and compared the optimal state values from Q learning to the optimal state values received from the DP. Each parameter was tested for 500 episodes.
   Tested Parameters:

   **Epsilon = 0.9**

   Constant Alpha = [0.01,0.1,0.2,**0.3**,0.4,0.5,0.6,0.7,0.8,0.9,1.0]

   Epsilon decay functions = [**1/((n/10)+1)**, 0.9*epsilon]

   The optimal solution (=best approximation) received by implementing the values appears in bold. After finding the optimal parameters we run the algorithm for 5000 episodes.

   The optimal solution for MC control described in figures bellow:

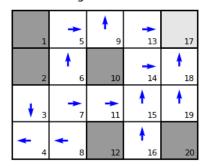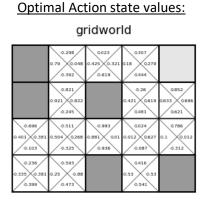Optimal state values:

Optimal Action state values:

Optimal policy:



we can see that this approximation is not converging to the optimal solution, and there for will suggest increasing number of episodes (CPU limitation)

MDP transition model for Action = North: (states notation from 0 to 19)

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0.8 | 0.1 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0.8 | 0.1 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0.1 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0.1 | 0 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0.8 | 0.1 | 0 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

MDP transition model for Action = East:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|----|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.8 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0.8 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0.8 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.8 | 0.1 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0.8 | 0.1 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

MDP transition model for Action = South:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0.1 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0.9 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.1 | 0.8 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 | 0 | 0 | 0 | 0.1 | 0.8 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

MDP transition model for Action = West:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0.1 | 0.8 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0.1 | 0.9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0.8 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 | 0 | 0 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0 | 0 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.8 | 0 | 0 | 0.1 | 0 | 0.1 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |