

Deterministic optimal policy for 2 Q-LEARNING after 100000 iterations.

