

Homework 5. Reinforcement learning

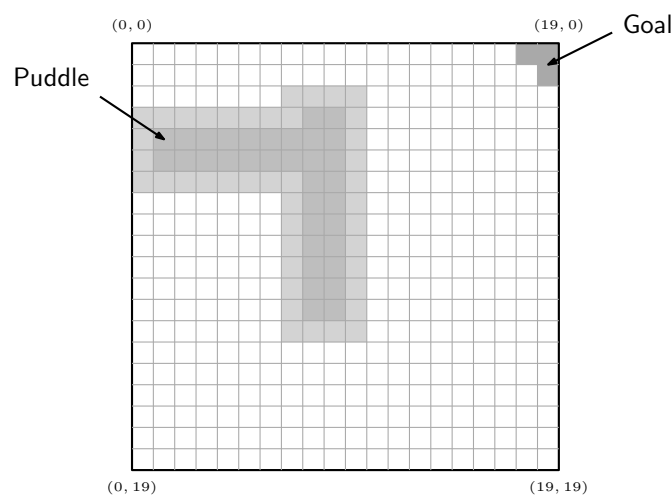


Figure 1: Puddle world.

Consider an all terrain vehicle navigating the grid world depicted in Fig. 1. The three shaded cells in the upper right corner correspond to the goal state, while the L-shaped shaded cells in the middle of the grid correspond to a puddle in which the vehicle may get stuck and damaged.

The vehicle has available the standard four actions—up, down, left and right. Each action

- Succeeds and moves the vehicle to the adjacent cell in the corresponding direction with a probability of 0.92;
- Fails and moves the vehicle to any of the other 3 adjacent cells with a probability of 0.02;
- Fails and the vehicle remains in the same cell with a probability of 0.02.

See Fig. 2 for an illustration of some movement situations.

Exercise 1.

- (a) Indicate the Q -values after a Q -learning update with step-size $\alpha = 0.1$, resulting from the transition at time step t .
- (b) Indicate the Q -values after a SARSA update with step-size $\alpha = 0.1$, resulting from the transition at time step t .
- (c) Explain the difference between on-policy and off-policy learning using Questions 1a and 1b to illustrate your explanation.