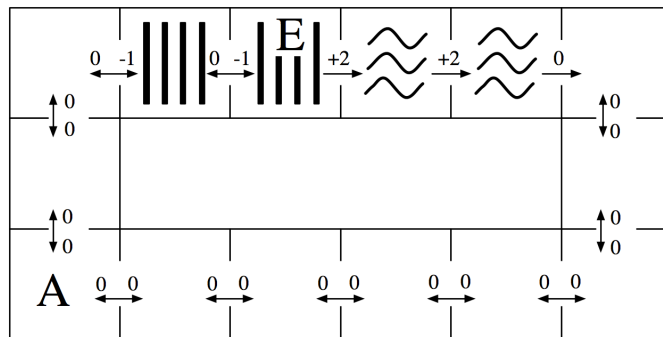# Section Handout 5

## Q1. MDPs: Grid-World Water Park

Consider the MDP drawn below. The state space consists of all squares in a grid-world water park. There is a single waterslide that is composed of two ladder squares and two slide squares (marked with vertical bars and squiggly lines respectively). An agent in this water park can move from any square to any neighboring square, unless the current square is a slide in which case it must move forward one square along the slide. The actions are denoted by arrows between squares on the map and all deterministically move the agent in the given direction. The agent cannot stand still: it must move on each time step. Rewards are also shown below: the agent feels great pleasure as it slides down the water slide (+2), a certain amount of discomfort as it climbs the rungs of the ladder (-1), and receives rewards of 0 otherwise. The time horizon is infinite; this MDP goes on forever.
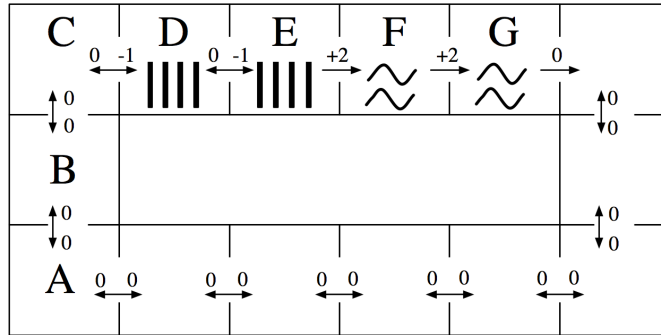


**(a)** How many (deterministic) policies $\pi$ are possible for this MDP?

**(b)** Fill in the blank cells of this table with values that are correct for the corresponding function, discount, and state. *Hint: You should not need to do substantial calculation here.*

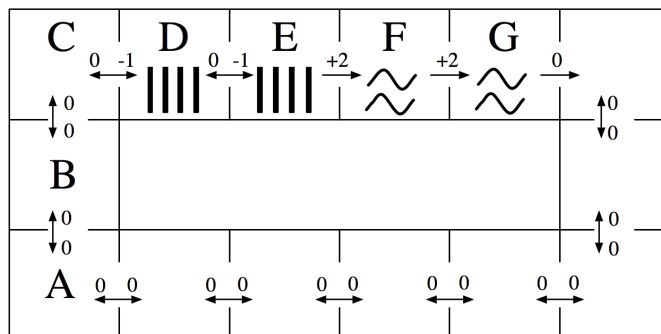|  | $\gamma$ | $s = A$ | $s = E$ |
|---|---|---|---|
| $V_3^*(s)$ | 1.0 | | |
| $V_{10}^*(s)$ | 1.0 | | |
| $V_{10}^*(s)$ | 0.1 | | |
| $Q_1^*(s, \text{west})$ | 1.0 | —— | |
| $Q_{10}^*(s, \text{west})$ | 1.0 | —— | |
| $V^*(s)$ | 1.0 | | |
| $V^*(s)$ | 0.1 | | |

Use this labeling of the state space to complete the remaining subproblems:



**(c)** Fill in the blank cells of this table with the Q-values that result from applying the Q-update for the transition specified on each row. You may leave Q-values that are unaffected by the current update blank. Use discount $\gamma = 1.0$ and learning rate $\alpha = 0.5$. Assume all Q-values are initialized to 0. (Note: the specified transitions would not arise from a single episode.)

| | | $Q(D, \text{west})$ | $Q(D, \text{east})$ | $Q(E, \text{west})$ | $Q(E, \text{east})$ |
|---|---|---|---|---|---|
| Initial: | | 0 | 0 | 0 | 0 |
| Transition 1: | $(s = D, a = \text{east}, r = -1, s' = E)$ | | | | |
| Transition 2: | $(s = E, a = \text{east}, r = +2, s' = F)$ | | | | |
| Transition 3: | $(s = E, a = \text{west}, r = 0, s' = D)$ | | | | |
| Transition 4: | $(s = D, a = \text{east}, r = -1, s' = E)$ | | | | |

The agent is still at the water park MDP, but now we're going to use function approximation to represent Q-values. Recall that a policy $\pi$ is *greedy* with respect to a set of Q-values as long as $\forall a, s\ Q(s, \pi(s)) \geq Q(s, a)$ (so ties may be broken in any way).
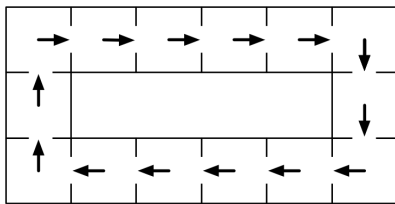
For the next subproblem, consider the following feature functions:

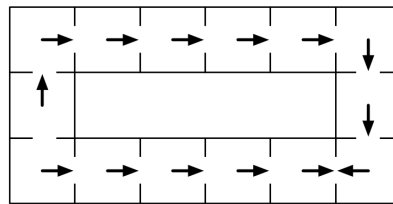$$f(s, a) = \begin{cases} 1 & \text{if } a = \text{east,} \\ 0 & \text{otherwise.} \end{cases}$$

$$f'(s, a) = \begin{cases} 1 & \text{if } (a = \text{east}) \wedge \text{isSlide}(s), \\ 0 & \text{otherwise.} \end{cases}$$

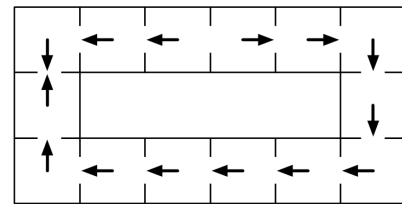(Note: isSlide($s$) is true iff the state $s$ is a slide square, i.e. either $F$ or $G$.)

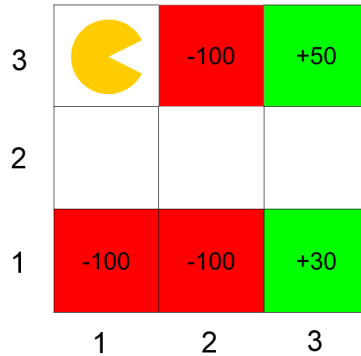Also consider the following policies:

$\pi_1$  $\pi_2$  $\pi_3$

(d) Which are greedy policies with respect to the Q-value approximation function obtained by running the single Q-update for the transition $(s = F, a = \text{east}, r = +2, s' = G)$ while using the specified feature function? You may assume that all feature weights are zero before the update. Use discount $\gamma = 1.0$ and learning rate $\alpha = 1.0$. Circle all that apply.

| $f$ | $\pi_1$ | $\pi_2$ | $\pi_3$ |
|---|---|---|---|
| $f'$ | $\pi_1$ | $\pi_2$ | $\pi_3$ |

# Q2. Direct Evaluation

Consider the grid-world given below and an agent who is trying to learn the optimal policy. Rewards are only awarded for taking the *Exit* action from one of the shaded states. Taking this action moves the agent to the Done state, and the MDP terminates. Assume $\gamma = 1$ and $\alpha = 0.5$ for all calculations. All equations need to explicitly mention $\gamma$ and $\alpha$ if necessary.



**(a)** The agent starts from the top left corner and you are given the following episodes from runs of the agent through this grid-world. Each line in an Episode is a tuple containing $(s, a, s', r)$.

| Episode 1 | Episode 2 | Episode 3 | Episode 4 | Episode 5 |
|---|---|---|---|---|
| (1,3), S, (1,2), 0 | (1,3), S, (1,2), 0 | (1,3), S, (1,2), 0 | (1,3), S, (1,2), 0 | (1,3), S, (1,2), 0 |
| (1,2), E, (2,2), 0 | (1,2), E, (2,2), 0 | (1,2), E, (2,2), 0 | (1,2), E, (2,2), 0 | (1,2), E, (2,2), 0 |
| (2,2), E, (3,2), 0 | (2,2), S, (2,1), 0 | (2,2), E, (3,2), 0 | (2,2), E, (3,2), 0 | (2,2), E, (3,2), 0 |
| (3,2), N, (3,3), 0 | (2,1), Exit, D, -100 | (3,2), S, (3,1), 0 | (3,2), N, (3,3), 0 | (3,2), S, (3,1), 0 |
| (3,3), Exit, D, +50 | | (3,1), Exit, D, +30 | (3,3), Exit, D, +50 | (3,1), Exit, D, +30 |

Fill in the following Q-values obtained from direct evaluation from the samples:

$Q((3,2), \text{N}) = $ _____     $Q((3,2), \text{S}) = $ _____     $Q((2,2), \text{E}) = $ _____