

Direct Q Evaluation A

Consider the grid-world given below and an agent (yellow) moving using these actions: N-North, W-West, E-East, S-South, and a special action D-Depart in terminal states (Exit). Rewards are only awarded for taking the *Exit* action from one of the terminal states (green and red). Assume discount factor  $\gamma = 1$  for all calculations.

|   |     |     |     |
|---|-----|-----|-----|
| 3 |     | -30 | 140 |
| 2 |     |     |     |
| 1 | -70 | -30 | 40  |
|   | 1   | 2   | 3   |

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, (2,3), 0     | (2,2), E, (2,1), 0     | (2,2), E, (2,1), 0     | (2,2), E, (2,3), 0     | (2,2), E, (3,2), 0                           |
| (2,3), D, (Exit,), -30 | (2,1), D, (Exit,), -30 | (2,1), D, (Exit,), -30 | (2,3), D, (Exit,), -30 | (3,2), N, (3,3), 0<br>(3,3), D, (Exit,), 140 |

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

$$Q((1,2), E) = \underline{\hspace{2cm}} \qquad Q((1,2), S) = \underline{\hspace{2cm}} \qquad Q((2,2), E) = \underline{\hspace{2cm}}$$