Amirabbas Rezasoltani      HwI      61039205

11, 13, 2023      AI

1)a) To find optimal policy in MDP, Q-learning is used. Suppose

in one of the episods, agent is in a particular state S. So

based on updating function of q-learning, $Q(s,a) = (1-\alpha)Q(s,a) +$

$\alpha(r_{s'} + \gamma \max\limits_{a'} Q(s',a'))$, which s' is next state after

doing action a. $r_{s'}$ is changed to $r_{s'} + c$. Assume all

cells of q-table are updated sufficient times. To update

$Q(s,a)$, $Q(s,a) = (1-\alpha)Q(s,a) + \alpha(r_{s'} + c + \gamma \max\limits_{a'} Q(s',a'))$.

assume $Q(s,a)$ is changed by $n_{s,a}$. So $Q(s,a)$ is changed

$\alpha c + \alpha\gamma n_{s',a'}$. Also $\max\limits_{a'} Q(s',a')$ is changed by, say

$\alpha c + \alpha\gamma n_{s'',a''}$ where s'' is next state of $\max Q(s',a')$,

and a'' is $\max\limits_{a''} Q(s'',a'')$. So $Q(s,a)$ is changed

$\alpha c + \alpha\gamma(\alpha c + \alpha\gamma n_{s'',a''}) = \alpha c + \alpha^2\gamma c + \alpha^2\gamma^2 n_{s'',a''}$. Continueing

computing $n_{s'',a''}$, we can show update for $Q(s,a)$ is

$\alpha c + \alpha^2\gamma c + \alpha^3\gamma^2 c + \dots = \alpha c(1 + \alpha\gamma + \alpha^2\gamma^2 + \dots) = \boxed{\alpha c \times \dfrac{1}{1-\alpha\gamma}}$.

since $Q(s,a)$ was arbitraly, all cells of table are added

with $\dfrac{\alpha c}{1-\alpha\gamma}$. If update continuess, this added value would be

Cumalative and will be a constant cofficient of $\frac{\alpha c}{1-\alpha \lambda}$. So optimal policy doesn't change.

b) Assume all cells of Q-table are uldated after multiplying all rewards by $c$. So to update $Q(s,a)$, we have,

$$Q(s,a) = (1-\alpha)Q(s,a) + \alpha(r_s \times c + \lambda \max_{a'} Q(s,a')).$$

assume $Q(s,a')$ is added $il$ by $n_{s,a'}$. So $Q(s,a)$ is

changed by $(\alpha r_s c - \alpha r_s) + \lambda \alpha n_{s',a'}$. So $Q(s,a)$ is where

changed by $\alpha r_s (c-1) + \lambda \alpha n_{s,a'}$. Also $n_{s,a'}$ is

$\alpha r_{s'}(c-1) + \lambda \alpha n_{s'',a''}$  ($s''$,$a''$ is introduced in Prelious Part). So $Q(s,a)$ is uldated by

$$\alpha r_s (c-1) + \lambda \alpha (\alpha r_{s'}(c-1) + \lambda \alpha n_{s'',a''}) = \alpha r_s (c-1) + \lambda \alpha^2 r_{s'}(c-1) +$$

$\lambda^2 \alpha^2 n_{s'',a''}$. By continueing computations, $Q(s,a)$ is added by $\alpha(c-1)(r_s + \lambda \alpha r_{s'} + ...)$. It depends on value of $r_s$ that a state can be mole valueable of less

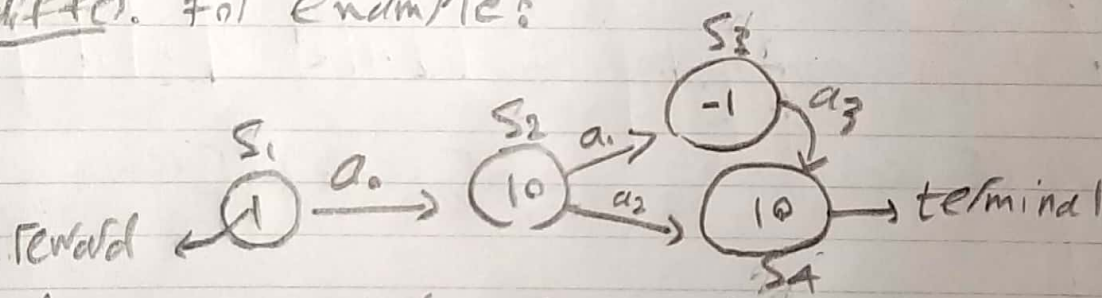valueable. If $\underline{c=1}$, optimal $\pi$ and values of $Q(s,a)$s don't differ. If $c \neq 1$, values and optimal Policy $\pi$ may differ.

If $\underline{c=0}$, all rewards become $0$ and so all cells of Q-table remain $0$. so all Policies are optimal.

c) Assume terminal state $S_{end}$. There is no valid action for that so $Q$s related to that won't update and its value is $0$. So like Part a, value change for $Q(s,a)$ is $dc(1+\alpha_2+\alpha_2 2^2+...+\alpha^k 2^k)$. Because steps to terminal state for each state differs, amount of update for $Q(s,a)$ differs. I mean if there are $k$ steps from state $s$ to terminal, update is $\alpha c(1+\alpha_2+...+\alpha^k 2^k)$. So optimal $\pi$ may differ. for example:



for diagram above, optimal Policy $\pi$ is:

$S_1 \rightarrow a_0$

$S_2 \rightarrow a_2$

$S_3 \rightarrow a_3$

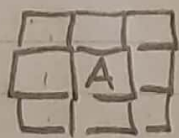But for $c=2$, optimal Policy $\pi$ will be $S_1 \rightarrow a_0$

$S_2 \rightarrow a_1$

$S_3 \rightarrow a_2$

2) a) To obtain what state we are in, check 8 cells around agent. In each cell, either there is wall, or food, or is empty or is invalid move. So number of states is minimum $4^8$. But some states are never happen. For example when all cells are invalid (out of bound), or all states are wall. so these states number can be reduced.

b) Actions : {up, right, bottom, left}. But in states which doing action causes going to invalid (out of bound) cells, this action is banned.

State : check 8 cells around agent like this.



where agent is in cell A. if cell is empty, its value is 0, if food is on that, value is 1, if wall is on that, its value is 2 and if it's out of bound, its value is 3. Using this, we can code all states.
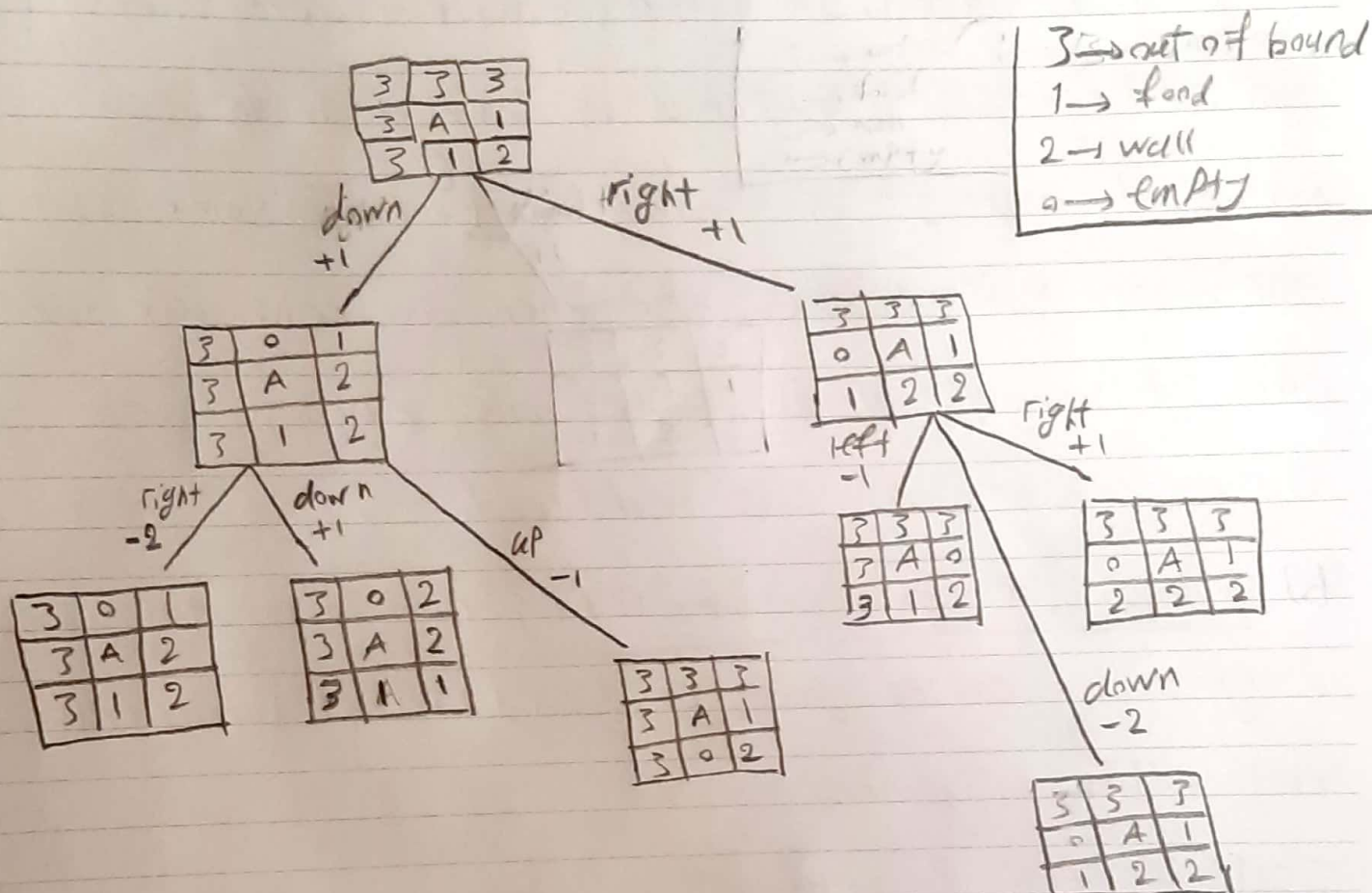
Reward : Based on value of each cell, reward function is as this: R(wall)=-2, R(food)=1, R(empty)=-1. The goal is find if all foods are eaten.

**3)** IS answered in code.

**4)** In episod 1, state of agent is like below for given map.

Legend:
- 3 → out of bound
- 1 → food
- 2 → wall
- 0 → empty



Tree diagram:

Root state:
| 3 | 3 | 3 |
| 3 | A | 1 |
| 3 | 1 | 2 |

down +1 →
| 3 | 0 | 1 |
| 3 | A | 2 |
| 3 | 1 | 2 |

right +1 →
| 3 | 3 | 3 |
| 0 | A | 1 |
| 1 | 2 | 2 |

From down state:

right -2 →
| 3 | 0 | 1 |
| 3 | A | 2 |
| 3 | 1 | 2 |

down +1 →
| 3 | 0 | 2 |
| 3 | A | 2 |
| 3 | 1 | 1 |

up -1 →
| 3 | 3 | 3 |
| 3 | A | 1 |
| 3 | 0 | 2 |

From right state:

left -1 →
| 3 | 3 | 3 |
| 3 | A | 0 |
| 3 | 1 | 2 |

right +1 →
| 3 | 3 | 3 |
| 0 | A | 1 |
| 2 | 2 | 2 |

down -2 →
| 3 | 3 | 3 |
| 0 | A | 1 |
| 1 | 2 | 2 |

**5)** As mentioned, to define states, 8 cells around agent
is checked and based on them, state is specified. So
there are $4^8$ states, since in each cell around, either
there is wall, or food, or is out of bound of is empty.
Also for each state, there are 4 moves (right, up, left, down).
But moving agent out of boundary is invalid so in q-table, $-\infty$ is
placed for action which cause invalid moves as symbol. After
updating q-table, for first state (root in tree in question 4)s

down, right, up, left have Values 1.108, 2.667, -∞, -∞

in order. For right node of root, values are -0.33, 2.67, -∞, -0.31

in order. For left node of root, Values are

2.67, -0.6, -0.47, -∞ in order. These results is

gathered from game with λ = 0.25 and α = 0.2, for
map A.

6) In code, there are 2 version of maps. map A is

given in question and map B is arbitrary.

7) Graphic Part is added to code.

8) A version of Pacman with ghost is implemented. what

differs from simple Pacman is, in Q-table, there ceen

be $8^8$ states. Because in every cell, there can be

ghost or not. But it can be reduced. Because ghost is

unique. So there are 2 total case. Either ghost is

in 8 cells or not. If ghost is, so there are 8 cells

which ghost can be. So there are only $9 \times 4^8$ states.