

## Homework 3

*Student Name: David Chen**Computing ID: dzc5ta*

- You will need to submit your solutions in PDF to UVA-Collab.

I collaborated with Timothy Han (txh7es)

1. a)

States: healthy, sick, and toothless

Actions: Eat candy or eat Vegetables

$$T(\text{healthy, vegetables, healthy}) = 1$$

$$T(\text{healthy, candy, healthy}) = 1/4$$

$$T(\text{healthy, candy, sick}) = 3/4$$

$$T(\text{sick, vegetables, healthy}) = 1/4$$

$$T(\text{sick, vegetables, sick}) = 3/4$$

$$T(\text{sick, candy, sick}) = 7/8$$

$$T(\text{sick, candy, toothless}) = 1/8$$

$$R(\text{candy}) = +10 \text{ happiness}$$

$$R(\text{vegetables}) = +4 \text{ happiness}$$

b)

$\pi_1$ : Annie always eats candy.

$$V(s) = \sum T(s, a, s')[R(s, a, s') + \gamma V^{\pi_1}(s')]$$

$$V^{\pi_1}(\text{healthy}) = T(\text{healthy, candy, healthy})[R(\text{healthy, candy, healthy}) + \gamma V^{\pi_1}(\text{healthy})] + T(\text{healthy, candy, sick})[R(\text{healthy, candy, sick}) + \gamma V^{\pi_1}(\text{sick})]$$

$$V^{\pi_1}(\text{healthy}) = 1/4 * (10 + \gamma V^{\pi_1}(\text{healthy})) + 3/4 * (10 + \gamma V^{\pi_1}(\text{sick}))$$

$$V^{\pi_1}(\text{healthy}) = 53.8899$$

$$V(\text{sick}) = T(\text{sick, candy, sick})[R(\text{sick, candy, sick}) + \gamma V^{\pi_1}(\text{sick})] + T(\text{sick, candy, toothless})[R(\text{sick, candy, toothless}) + \gamma V^{\pi_1}(\text{toothless})]$$

$$V(\text{sick}) = 7/8 * (10 + \gamma V^{\pi_1}(\text{sick})) + 1/8 * (10 + \gamma V^{\pi_1}(\text{toothless}))$$

$$V^{\pi_1}(\text{sick}) = 47.0588$$

$$V(\text{toothless}) = 0$$

$\pi_2$ : Annie always eats vegetables.

$$V(s) = \sum T(s, a, s')[R(s, a, s') + \gamma V^{\pi_2}(s')]$$

$$V^{\pi_2}(\text{healthy}) = T(\text{healthy, vegetables, healthy})[R(\text{healthy, vegetables, healthy}) + \gamma V^{\pi_2}(\text{healthy})]$$

$$V^{\pi_2}(\text{healthy}) = 1 * (4 + \gamma V^{\pi_2}(\text{healthy}))$$

$$V^{\pi_2}(\text{healthy}) = 40$$

$$V^{\pi_2}(\text{sick}) = T(\text{sick}, \text{vegetables}, \text{healthy})[R(\text{sick}, \text{vegetables}, \text{healthy}) + \gamma V^{\pi_2}(\text{sick})] +$$

$$T(\text{sick}, \text{vegetables}, \text{sick})[R(\text{sick}, \text{vegetables}, \text{sick}) + \gamma V^{\pi_2}(\text{toothless})]$$

$$V^{\pi_2}(\text{sick}) = 1/4 * (4 + \gamma V^{\pi_2}(\text{healthy})) + 3/4 * (4 + \gamma V^{\pi_2}(\text{sick}))$$

$$V^{\pi_2}(\text{sick}) = 40$$

c)

$\pi_1$ : Annie always eats candy.

$$V^{\pi_1}(\text{toothless}) = 0$$

$$V^{\pi_1}(\text{sick}) = 47.0588$$

$$V^{\pi_1}(\text{healthy}) = 53.8899$$

$$\pi^*(\text{sick}): \text{candy: } 7/8 * (10 + .9 * 47.0588) + 1/8 * (10 + .9 * 0) = 47.058805$$

$$\text{vegetable: } 3/4 * (4 + .9 * 47.0588) + 1/4 * (4 + .9 * 53.8899) = 47.8899175$$

$$\pi^*(\text{sick}) = \text{eat vegetable}$$

$$\pi^*(\text{healthy}): \text{candy: } 3/4 * (10 + .9 * 47.0588) + 1/4 * (10 + .9 * 53.8899) =$$

$$53.8899175$$

$$\text{vegetable: } 1 * (4 + .9 * 53.8899) = 52.50091$$

$$\pi^*(\text{sick}) = \text{eat candy}$$

$\pi_2$ : Annie eats vegetables when she is sick and candy when she is healthy.

$$V^{\pi_2}(\text{toothless}) = 0$$

$$V^{\pi_2}(\text{sick}) = 51.1685$$

$$V^{\pi_2}(\text{healthy}) = 57.1685$$

$$\pi^*(\text{sick}): \text{candy: } 7/8 * (10 + .9 * 51.1685) + 1/8 * (10 + .9 * 0) = 50.29519375$$

$$\text{vegetable: } 3/4 * (4 + .9 * 51.1685) + 1/4 * (4 + .9 * 57.1685) = 51.40165$$

$$\pi^*(\text{sick}) = \text{eat vegetable}$$

$$\pi^*(\text{healthy}): \text{candy: } 3/4 * (10 + .9 * 51.1685) + 1/4 * (10 + .9 * 57.1685) = 57.40165$$

$$\text{vegetable: } 1 * (4 + .9 * 57.1685) = 55.45165$$

$$\pi^*(\text{sick}) = \text{eat candy}$$

$\pi_3$ : Annie eats vegetables when she is sick and candy when she is healthy.

The optimal policy has been found since it stayed the same. Therefore, the policy after the third iteration will still be: Annie eats vegetables when she is sick and candy when she is healthy.

d)

i) Policy converges long before values converge.

iii) There must be a finite number of policies, since there is a finite number of states and actions.

v) If there are two actions that lead to the same state than either action can be in an optimal policy.

2. a)

$$V^\pi(s') = V^\pi(s') + \alpha(R(s', \pi(s'), s'') + \gamma V^\pi(s''))$$

$$V^\pi(s) = V^\pi(s) + \alpha(R(s, \pi(s), s') + \gamma(V^\pi(s') + \alpha(R(s', \pi(s'), s'') + \gamma V^\pi(s''))))$$

b)

i) The first four (state, action) pairs would be:

$$(S3, A3) \rightarrow (S5, A1) \rightarrow (S6, A3) \rightarrow (S5, A1)$$

ii) If only the best Q-value is chosen, there might be unexplored regions. We can alter this by sometimes, choosing a sub-optimal option. This way, other regions can be explored.

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$$

$$Q^*(S3, A1) = 1 * 16 = 16$$

$$Q^*(S2, A2) = 0.5 * 4 + 0.5 * 12 = 8$$

$$Q^*(S2, A1) = 1 * 6 = 6$$

$$Q^*(S1, A2) = 0.25 * 8 + 0.75 * 16 = 14$$

$$Q^*(S1, A1) = 0.5 * 8 + 0.5 * 16 = 12$$

The optimal policy: At state S1 and S2, take action A2 since that action maximizes the value at those states.

$$Q = 0$$

$$Q \leftarrow (1 - 0.5)(0) + 0.5 * 12 = 6$$

$$Q \leftarrow (1 - 0.5)(6) + 0.5 * 16 = 11$$

$$Q \leftarrow (1 - 0.5)(11) + 0.5 * 6 = 8.5$$

3. In beam search, there are k beams, or best paths, at the end of each iteration. For each iteration, each beam is expanded and different paths are explored. Each path is then ranked and the best k paths become the new starting beams. This process repeats until a goal is reached. Beam search can model the ant colony optimization problem because each beam can be a path that is being explored. The ants will when leave more pheromone based on how much potential the food on that path has. The amount of pheromone on a path can be used to rank each path. At the end of the day ants return to their nest after a long day of gathering food. At the beginning of the next day, the ants start exploring from the top k paths with the most pheromone.