

## Stanford CME 241 (Winter 2021) - Assignment 2

### Assignments:

1. For a *Deterministic* Policy  $\pi_D : \mathcal{S} \rightarrow \mathcal{A}$ , write with precise mathematical notation the 4 MDP Bellman Policy Equations, i.e.,  $V^{\pi_D}$  in terms of  $Q^{\pi_D}$ ,  $Q^{\pi_D}$  in terms of  $V^{\pi_D}$ ,  $V^{\pi_D}$  in terms of  $Q^{\pi_D}$ ,  $Q^{\pi_D}$  in terms of  $V^{\pi_D}$ . Note that in the book, we have written the 4 MDP Policy Equations in terms of the notation for a stochastic policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ .
2. Consider an MDP with an infinite set of states  $\mathcal{S} = \{1, 2, 3, \dots\}$ . The start state is  $s = 1$ . Each state  $s$  allows a continuous set of actions  $a \in [0, 1]$ . The transition probabilities are given by:

$$\mathbb{P}[s + 1 \mid s, a] = a, \mathbb{P}[s \mid s, a] = 1 - a \text{ for all } s \in \mathcal{S} \text{ for all } a \in [0, 1]$$

For all states  $s \in \mathcal{S}$  and actions  $a \in [0, 1]$ , transitioning from  $s$  to  $s + 1$  results in a reward of  $1 - a$  and transitioning from  $s$  to  $s$  results in a reward of  $1 + a$ . The discount factor  $\gamma = 0.5$ .

- Calculate the Optimal Value Function  $V^*(s)$  for all  $s \in \mathcal{S}$
  - Calculate an Optimal Deterministic Policy  $\pi^*(s)$  for all  $s \in \mathcal{S}$
3. Consider an array of  $n + 1$  lilypads on a pond, numbered 0 to  $n$ . A frog sits on a lilypad other than the lilypads numbered 0 or  $n$ . When on lilypad  $i$  ( $1 \leq i \leq n - 1$ ), the frog can croak one of two sounds  $A$  or  $B$ . If it croaks  $A$  when on lilypad  $i$  ( $1 \leq i \leq n - 1$ ), it is thrown to lilypad  $i - 1$  with probability  $\frac{i}{n}$  and is thrown to lilypad  $i + 1$  with probability  $\frac{n-i}{n}$ . If it croaks  $B$  when on lilypad  $i$  ( $1 \leq i \leq n - 1$ ), it is thrown to one of the lilypads  $0, \dots, i - 1, i + 1, \dots, n$  with uniform probability  $\frac{1}{n}$ . A snake, perched on lilypad 0, will eat the frog if the frog lands on lilypad 0. The frog can escape the pond (and hence, escape the snake!) if it lands on lilypad  $n$ .

What should the frog croak when on each of the lilypads  $1, 2, \dots, n - 1$ , in order to maximize the probability of escaping the pond (i.e., reaching lilypad  $n$  before reaching lilypad 0)? Although there are more than one ways of solving this problem, we'd like to solve it by modeling it as an MDP and identifying the Optimal Policy.

- Express with clear mathematical notation the state space, action space, transitions function and rewards function of an MDP so that the above *frog-escape* problem would be solved by arriving at the Optimal Value Function (and hence, the Optimal Policy) of this MDP.
- Write code to model this MDP as an instance of the `FiniteMarkovDecisionProcess` class. We have learnt that there exists an optimal deterministic policy, and there are  $2^n$  possible deterministic policies for this problem. Write code to create each of these  $2^n$  deterministic policies (as instances of `FinitePolicy` class), create a policy-implied Finite MRP for each of these deterministic policies (using the `apply_finite_policy` method of `FiniteMarkovDecisionProcess` class), and evaluate the Value Function for each of those implied Finite MRPs (using the `get_value_function_vec` method of `FiniteMarkovRewardProcess` class). This should give you the Optimal Value Function and the Optimal Deterministic Policy.
- Using your code, plot a graph of the Optimal Escape-Probability and of the associated Optimal Croak, as a function of the states of this MDP, for  $n = 3, n = 6$  and  $n = 9$ . By looking at the results on this graph, what pattern do you observe for the optimal policy as you vary  $n$  from 3 to 9?

4. Consider a continuous-states, continuous-actions, discrete-time, non-terminating MDP with state space as  $\mathbb{R}$  and action space as  $\mathbb{R}$ . When in state  $s \in \mathbb{R}$ , upon taking action  $a \in \mathbb{R}$ , one transitions to next state  $s' \in \mathbb{R}$  according to a normal distribution  $s' \sim \mathcal{N}(s, \sigma^2)$  for a fixed variance  $\sigma^2 \in \mathbb{R}^+$ . The corresponding cost associated with this transition is  $e^{as'}$ , i.e., the cost depends on the action  $a$  and the state  $s'$  one transitions to. The problem is to minimize the infinite-horizon *Expected Discounted-Sum of Costs* (with discount factor  $\gamma < 1$ ). For this assignment, solve this problem just for the special case of  $\gamma = 0$  (i.e., the myopic case) using elementary calculus. Derive an analytic expression for the optimal action in any state and the corresponding optimal cost.

**Optional:** How would you approach this problem for the general case of  $0 \leq \gamma < 1$ ?