

## CS533

### Intelligent Agents and Decision Making MDP Basics and Finite Horizon Problems

1. Construct a simple Markov Decision Process such that the optimal policy for maximizing finite-horizon total reward must be non-stationary. That is, the MDP should not have a stationary policy that maximizes the finite horizon total reward.
2. In many problems, not all actions are applicable in all states and many actions only lead to a small number of next states, compared to the total number of states. In this question, we consider how the complexity of finite-horizon value iteration and policy evaluation can be improved for such problems.

To capture the notion of applicable actions, suppose that we have a function  $\text{LEGAL}(s)$  that takes a state  $s$  and returns the set of legal actions in  $s$ . Also suppose that we have a function  $\text{NEXT}(s, a)$ , which takes a state  $s$  and action  $a$  as input and returns the set of states that have non-zero probability of occurring after taking  $a$  in state  $s$ . That is,

$$\text{NEXT}(s, a) = \{s' \mid T(s, a, s') > 0\}.$$

Assume that we are considering an MDP with  $n$  states and  $m$  actions such that for any state  $s$  and action  $a$  we have  $|\text{LEGAL}(s)| \leq k$  and  $|\text{NEXT}(s, a)| \leq r$ . Assume that the time and space complexity of evaluating the functions  $\text{NEXT}$  and  $\text{LEGAL}$  are linear in the sizes of their output (i.e. the number of elements in their sets).

- (a) Describe how to modify the finite-horizon policy evaluation algorithm described in class, using one or both of the new functions, so that the time complexity is improved when  $r < n$  and  $k < m$ . What is the time complexity? The time complexity should be expressed in terms of  $r$  and  $k$  when possible and may also involve  $n$  and  $m$ .
  - (b) Repeat part (a) but for the finite-horizon value iteration algorithm described in class.
3. Our basic definition of an MDP in class defined the reward function  $R(s)$  to be a function of just the state, which we will call a *state reward function*. It is also common to define a reward function to be a function of the state and action, written as  $R(s, a)$ , which we will call a *state-action reward function*. The meaning is that the agent gets a reward of  $R(s, a)$  when they take action  $a$  in state  $s$ . While this may seem to be a significant difference, it does not fundamentally extend our modeling power, nor does it fundamentally change the algorithms that we have developed.
    - (a) Describe a real world problem where the corresponding MDP is more naturally modeled using a state-action reward function compared to using a state reward function.
    - (b) Modify the finite-horizon value iteration algorithm so that it works for state-action reward functions. Do this by writing out the new update equation that is used each iteration and explaining the modification from the equation given in class for state rewards.
    - (c) Any MDP with a state-action reward function can be transformed into an “equivalent” MDP with just a state reward function. Show how any MDP with a state-action reward function  $R(s, a)$  can be transformed into a different MDP with state reward function  $R(s)$ , such that the optimal policies in the new MDP correspond exactly to the optimal policies in the original MDP. That is an optimal policy in the new MDP can be mapped

to an optimal policy in the original MDP. *Hint: It will be necessary for the new MDP to introduce new “book keeping” states that are not in the original MDP.*

4. (***k*-th order MDPs.**) A standard MDP is described by a set of states  $S$ , a set of actions  $A$ , a transition function  $T$ , and a reward function  $R$ . Where  $T(s, a, s')$  gives the probability of transitioning to  $s'$  after taking action  $a$  in state  $s$ , and  $R(s)$  gives the immediate reward of being in state  $s$ .

A  $k$ -order MDP is described in the same way with one exception. The transition function  $T$  depends on the current state  $s$  and also the previous  $k-1$  states. That is,  $T(s_{k-1}, \dots, s_1, s, a, s') = \Pr(s'|a, s, s_1, \dots, s_{k-1})$  gives the probability of transitioning to state  $s'$  given that action  $a$  was taken in state  $s$  and the previous  $k-1$  states were  $(s_{k-1}, \dots, s_1)$ .

Given a  $k$ -order MDP  $M = (S, A, T, R)$  describe how to construct a standard (first-order) MDP  $M' = (S', A', T', R')$  that is equivalent to  $M$ . Here equivalent means that a solution to  $M'$  can be easily converted into a solution to  $M$ . Be sure to describe  $S'$ ,  $A'$ ,  $T'$ , and  $R'$ . Give a brief justification for your construction.

5. Suppose that in a finite-horizon setting, we would like the reward function to depend on the time-to-go. That is, the reward function will be of the form  $R(s, t)$ , which says that we get reward  $R(s, t)$  for being in state  $s$  when the time-to-go is  $t$ . Can finite-horizon value iteration be modified to take this reward function into account? If so, show how to modify the equations. If not, then give an argument why.