

**Exam Markov Decision Theory  
and Algorithmic Methods (191531920)**

April 7, 2014    8:45-11:45 hrs

This exam consists of 4 exercises.  
Motivate all your answers.

1. Consider an infinite horizon average reward Markov Decision Problem (MDP) with state space  $S = \{s_1, s_2\}$ , and action sets  $A_{s_1} = \{a_{1,1}, a_{1,2}\}$ ,  $A_{s_2} = \{a_{2,1}, a_{2,2}\}$ . The immediate rewards are  $r(s_1, a_{1,1}) = 4$ ,  $r(s_1, a_{1,2}) = 6$ ,  $r(s_2, a_{2,1}) = -4$ ,  $r(s_2, a_{2,2}) = -6$ . The transition probabilities are given by  $p(s_2|s_1, a_{1,1}) = 1/2$ ,  $p(s_2|s_1, a_{1,2}) = 1$ ,  $p(s_2|s_2, a_{2,1}) = 1/2$ , and  $p(s_2|s_2, a_{2,2}) = 0$ .

- (a) The stationary policy  $d^\infty$  is defined by the decision rule  $d$  satisfying  $d(s_1) = a_{1,2}$  and  $d(s_2) = a_{2,2}$ . Calculate the gain  $g^{d^\infty}$  of this stationary policy.
- (b) The optimality equations in vector notation are given by  $B(g, h) = 0$  with

$$B(g, h)(s) = \max_{a \in A_s} \left\{ r(s, a) - g + \sum_{j \in S} p(j|s, a)h(j) - h(s) \right\}.$$

Write down the optimality equations for this MDP. Use these equations and their properties to show that the optimal gain  $g^*$  is bounded, namely  $-4 \leq g^*(s) \leq 6$ .

- (c) Show that  $g^*(s) = 2/3$ ,  $s \in S$ , is the optimal gain.
- (d) Let  $h = (10/3, -10/3)$  be a bias vector. Determine a decision rule  $d$  that is  $h$ -improving.
- (e) Suppose that you are asked to check whether or not a given policy  $\pi$  is average optimal. Mention two different ways to do so.
2. (a) Consider an infinite horizon discounted MDP. Explain in words why we may restrict attention to Markov policies, instead of history-dependent policies, when analyzing discounted MDPs.
- (b) One algorithm for solving infinite horizon discounted MDPs is the value iteration algorithm. This results in a stationary policy  $(d_\varepsilon)^\infty$  with

$$d_\varepsilon(s) \in \arg \max_{a \in A_s} \left\{ r(s, a) + \sum_{j \in S} \lambda p(j|s, a)v^{n+1}(j) \right\}$$

for each state  $s \in S$ . Prove that the policy  $(d_\varepsilon)^\infty$  is  $\varepsilon$ -optimal.

3. Consider large-scale MDPs with countable state space  $S = \{0, 1, \dots\}$ , discount factor  $\lambda$ , and unbounded rewards.

(a) In this setting, the weighted supremum norm with respect to  $w$  is used:  $\|v\|_w = \sup_{s \in S} w(s)^{-1}|v(s)|$  with  $w$  an arbitrary positive real-valued function on  $S$  satisfying  $\inf_{s \in S} w(s) > 0$ .

Let  $A_s = \{0, 1, 2, \dots, M\}$  for all states  $s$ ,  $r(s, a) = s$ , and  $p(j|s, a) = 1$  if  $j = s + a$  and  $p(j|s, a) = 0$  else. Let  $w(s) = \max(s, 1)$ . Show that there exists a constant  $\kappa$ ,  $0 \leq \kappa < \infty$ , such that

$$\sum_{j \in S} p(j|s, a)w(j) \leq \kappa w(s), \quad \text{for all } a \in A_s, \text{ for all } s \in S.$$

(This is one of the conditions for existence of an optimal policy.)

(b) Under suitable conditions (one of them is mentioned in part (a)), the optimality equation has an optimal solution; that is, the MDP has a value. Why do algorithms like the value iteration algorithm not work in this case? And, how can we approximate the value in practice?

4. Approximate dynamic programming may be applied to average cost problems. One method for solving such problems is the approximate policy evaluation, for approximating the cost of a stationary policy  $\mu$ . Let  $x_k$  denote the state at decision epoch  $k$ ,  $p_{ij}$  the transition probability of going from state  $i$  to state  $j$  given the policy  $\mu$ , and  $g(x_k, x_{k+1})$  the immediate cost at decision epoch  $k$  starting state  $x_k$ , using policy  $\mu$  and the next state is  $x_{k+1}$ .

The optimality equations are

$$h(i) = \sum_{j=1}^n p_{ij} (g(i, j) - \eta + h(j)),$$

with  $\eta$  the average cost for each initial state. Assume that  $\eta$  is known. The goal is to approximate the vector  $h$  by a linear architecture  $\tilde{h}(i, r) = \phi(i)'r$ .

(a) Explain what a linear architecture is.

(b) Explain how to use a projected equation to approximate  $h$  by a linear architecture.

Points:

1					2		3		4		Total
a	b	c	d	e	a	b	a	b	a	b	
2	3	4	4	2	2	4	4	4	3	4	+ 4 = 40