

Please use the L<sup>A</sup>T<sub>E</sub>X template to produce your writeups. See the Homework Assignments page on the class website for details. Hand in at: <https://webhandin.eng.utah.edu/index.php>.

## 1 Policy Iteration (30pts)

The card game of high-low is played with an infinite deck whose only cards are 2, 3, and 4 in equal proportion. You start with one of the cards showing, and say either *high* or *low*. Then a new card is flipped, and you compare the value of the new card to that of the old card.

- If you are right, you get the value of the new card.
- If the new card has the same value, you don't get any points.
- If you are wrong, the game is done.

If you are not done, the new card then becomes the reference card for drawing the next new card. You accumulate points as above until you are wrong and the game ends.

1. Formulate high-low as an MDP, by listing the states, actions, transition rewards, and transition probabilities.

$s$	$a$	$s'$	$T(s, a, s')$
2	<i>high</i>	<i>win</i>	2/3
2	<i>high</i>	<i>equal</i>	1/3
2	<i>low</i>	<i>lose</i>	2/3
2	<i>low</i>	<i>equal</i>	1/3
3	<i>high</i>	<i>win</i>	1/3
3	<i>high</i>	<i>equal</i>	1/3
3	<i>high</i>	<i>lose</i>	1/3
3	<i>low</i>	<i>win</i>	1/3
3	<i>low</i>	<i>equal</i>	1/3
3	<i>low</i>	<i>lose</i>	1/3
4	<i>high</i>	<i>lose</i>	2/3
4	<i>high</i>	<i>equal</i>	1/3
4	<i>low</i>	<i>win</i>	2/3
4	<i>low</i>	<i>equal</i>	1/3

Transition Model

$s'$	$R(s')$
<i>win</i>	{2, 3, 4}
<i>equal</i>	0
<i>lose</i>	terminate

Rewards

2. You will be doing one iteration of policy iteration. Assume the initial policy  $\pi_0(s) = \text{high}$ .

- (a) Perform policy evaluation to solve for the utility values  $V^{\pi_0}(s)$  for the appropriate states  $s$ . Please solve these equations analytically.

$$V^{\pi_0}(s) = \sum_{s'} T(s, \pi_0, s') [R(s, \pi_0(s), s') + \gamma V^{\pi_0}(s')]$$

We can also calculate the expected value of drawing a new card by taking the average of the values since they're all equally likely. We assume  $\gamma = 1$

$$V_{k+1}^{\pi_0}(2) = \frac{1}{3} (0 + V_k^{\pi_0}(2)) + \frac{1}{3} (3 + V_k^{\pi_0}(3)) + \frac{1}{3} (4 + V_k^{\pi_0}(4))$$

But we don't know the values for  $V_k^{\pi_0}(3)$  or  $V_k^{\pi_0}(4)$ , so we can solve those first and then plug them back into the above equation. We get

$$\begin{aligned} V_k^{\pi_0}(3) &= \frac{1}{3}(\text{terminate}) + \frac{1}{3} (0 + V_k^{\pi_0}(3)) + \frac{1}{3} (4 + V_k^{\pi_0}(4)) \\ V_k^{\pi_0}(4) &= \frac{2}{3}(\text{terminate}) + \frac{1}{3} (0 + V_k^{\pi_0}(4)) \end{aligned}$$

Where we can ignore the termination states as we want to *maximize* our output. We can solve for  $V_k^{\pi_0}(4)$  and then  $V_k^{\pi_0}(3)$  to give

$$\begin{aligned} V_k^{\pi_0}(4) &= 0 \\ V_k^{\pi_0}(3) &= \frac{1}{3} (0 + V_k^{\pi_0}(3)) + \frac{1}{3} (4 + V_k^{\pi_0}(4)) \\ &= \frac{1}{3} V_k^{\pi_0}(3) + \frac{4}{3} \\ &= 2 \end{aligned}$$

Finally, by solving for  $V_k^{\pi_0}(2)$

$$\begin{aligned} V_{k+1}^{\pi_0}(2) &= \frac{1}{3} (0 + V_k^{\pi_0}(2)) + \frac{1}{3} (3 + V_k^{\pi_0}(3)) + \frac{1}{3} (4 + V_k^{\pi_0}(4)) \\ &= \frac{9}{2} \end{aligned}$$

- (b) Perform policy improvement to find the next policy  $\pi_1(s)$ .

The equation for calculating the next policy is

$$\pi_{k+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^{\pi_k}(s')]$$

We can use this to find the next policy update for each state, giving

$$\pi_1(2) = \operatorname{argmax}_a \begin{cases} \frac{1}{3} \left(0 + \frac{9}{2}\right) + \frac{1}{3} \left(0 + \frac{9}{2}\right) & = 3/2 & (low) \\ \frac{1}{3} \left(0 + \frac{9}{2}\right) + \frac{1}{3} (3 + 2) + \frac{1}{3} (4 + 0) & = 27/6 & (high) \end{cases}$$

$$\pi_1(2) = high$$

$$\pi_1(3) = \operatorname{argmax}_a \begin{cases} \frac{1}{3} (0 + 2) + \frac{1}{3} \left(2 + \frac{9}{2}\right) & = 17/6 & (low) \\ \frac{1}{3} (0 + 2) + \frac{1}{3} (4 + 0) & = 2 & (high) \end{cases}$$

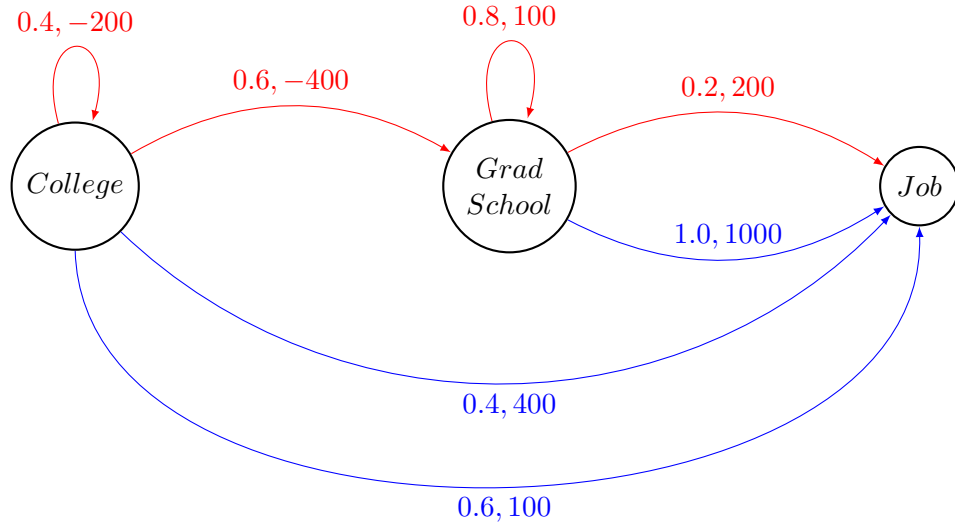
$$\pi_1(3) = low$$

$$\pi_1(4) = \operatorname{argmax}_a \begin{cases} \frac{1}{3} \left(2 + \frac{9}{2}\right) + \frac{1}{3} (3 + 2) + \frac{1}{3} (0 + 0) & = 23/6 & (low) \\ \frac{1}{3} (0 + 0) & = 0 & (high) \end{cases}$$

$$\pi_1(4) = high$$

## 2 Life as a Student

Jennifer is currently finishing up in college and is trying to decide what she wants to do with the rest of her life. She sketches her options as a known MDP:



At each point in her life she can choose to either continue in school (action  $x$ ) or try to get a job (action  $y$ ). Her three states are **College**, **Grad School** and **Job**. **J** is a terminal state. Transition probabilities  $t$  and immediate rewards  $r$  are shown on the arcs as  $t, r$ . The discount  $\gamma = 0.5$ .

1. Using value iteration, compute values for each state. Add or remove extra lines for iteration as needed.

$i$	$V_i^*(C)$	$V_i^*(G)$
0	0.0	0.0
1	220.0	1000.0
2	220.0	1000.0

The non-linear equation for *Value Iteration* is defined as

$$V_{k+1}^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')]$$

We can calculate the value of  $V$  for two iterations of *College*.

$$\begin{aligned}
 V_1^*(C) &= \max_a \begin{cases} 0.4(-200 + 0) + 0.6(-400 + 0) &= -320 & (x) \\ 0.4(400 + 0) + 0.6(100 + 0) &= 220 & (y) \end{cases} \\
 V_2^*(C) &= \max_a \begin{cases} 0.4(-200 + 0.5(220)) + 0.6(-400 + 0.5(220)) &= -210 & (x) \\ 0.4(400 + 0.5(0)) + 0.6(100 + 0.5(0)) &= 220 & (y) \end{cases}
 \end{aligned}$$

2. Suppose Jennifer didn't actually know the values of the MDP. Instead, she watched three of her older siblings go through life. They exhibited the following trajectories:

Sibling 1	Sibling 2	Sibling 3
$C, \textcolor{red}{x}, -200$	$C, \textcolor{red}{x}, -400$	$C, \textcolor{red}{x}, -400$
$C, \textcolor{red}{x}, -200$	$G, \textcolor{red}{x}, 100$	$G, \textcolor{red}{x}, 100$
$C, \textcolor{blue}{y}, 400$	$G, \textcolor{red}{x}, 100$	$G, \textcolor{red}{x}, 100$
$J$	$G, \textcolor{blue}{y}, 1000$	$G, \textcolor{red}{x}, 200$
	$J$	$J$

- (a) According to direct estimation, what are the transition probabilities  $T(C, \textcolor{red}{x}, C)$ ,  $T(C, \textcolor{red}{x}, G)$  and  $T(C, \textcolor{blue}{y}, J)$ ?

$T(C, \textcolor{red}{x}, C)$	$T(C, \textcolor{red}{x}, G)$	$T(C, \textcolor{blue}{y}, J)$
4/9	5/9	1/1

- (b) According to direct estimation, what are the rewards  $R(C, \textcolor{red}{x}, C)$ ,  $R(C, \textcolor{red}{x}, G)$  and  $R(C, \textcolor{blue}{y}, J)$ ?

$R(C, \textcolor{red}{x}, C)$	$R(C, \textcolor{red}{x}, G)$	$R(C, \textcolor{blue}{y}, J)$
-300	120	400

- (c) Use TD Learning instead to find estimates of the value function  $V^\pi(s)$ , assuming  $\alpha = 0.5$ .

The equation for Temporal Difference Learning is

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi, s') + \gamma V^\pi(s') - V^\pi(s)]$$