

Please use the L^AT_EX 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 Value Iteration

At the AI casino, there are two things to do: Eat Buffet and Play AI Blackjack. You start out Poor and Hungry, and would like to leave the casino Rich and Full. If you Play while you are Full you are more likely to become Rich, but if you are Poor you may have a hard time becoming Full on your budget. We can model your decision making process as the following MDP:

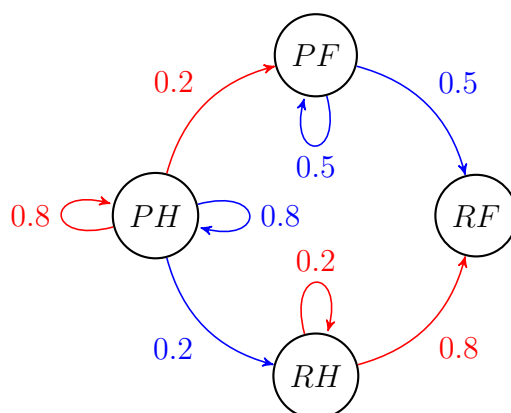
State Space {PoorHungry, PoorFull, RichHungry, RichFull}
 Actions {Eat, Play}
 Initial State PoorHungry
 Terminal State RichFull

s	a	s'	$T(s, a, s')$
PoorHungry	Play	PoorHungry	0.8
PoorHungry	Play	RichHungry	0.2
PoorHungry	Eat	PoorHungry	0.8
PoorHungry	Eat	PoorFull	0.2
PoorFull	Play	PoorFull	0.5
PoorFull	Play	RichFull	0.5
RichHungry	Eat	RichHungry	0.2
RichHungry	Eat	RichFull	0.8

Transition Model

s'	$R(s')$
PoorHungry	-1
PoorFull	1
RichHungry	0
RichFull	5

Rewards



Where **red** denotes the action to **Eat** and **blue** denotes **Play**.

1. Complete the table for the first 3 iterations of Value Iteration. Assume $\gamma = 1$.

State	$i = 0$	$i = 1$	$i = 2$	$i = 3$
PoorHungry	0.00	-1.00	-1.60	-1.98
PoorFull	0.00	1.00	1.50	1.75
RichHungry	0.00	0.00	0.00	0.00
RichFull	0.00	0.00	0.00	0.00

The values of *PoorHungry* (PH) is calculated for two iterations to show the steps.

$$V_1(PH) = R(PH) + \gamma \max_a \left\{ \begin{array}{l} T(PH, play, PH)V_0(PH) + T(PH, play, RH)V_0(RH) \\ T(PH, eat, PH)V_0(PH) + T(PH, eat, PF)V_0(PF) \end{array} \right.$$

$$V_1(PH) = -1.0 + (\gamma = 1) \max_a \left\{ \begin{array}{l} 0.8(-1) + 0.2(0) \\ 0.8(-1) + 0.2(0) \end{array} \right.$$

$$V_1(PH) = -1.0$$

$$V_2(PH) = R(PH) + \gamma \max_a \left\{ \begin{array}{l} T(PH, play, PH)V_1(PH) + T(PH, play, RH)V_1(RH) \\ T(PH, eat, PH)V_1(PH) + T(PH, eat, PF)V_1(PF) \end{array} \right.$$

$$V_2(PH) = -1.0 + (\gamma = 1) \max_a \left\{ \begin{array}{l} 0.8(-1) + 0.2(0) \\ 0.8(-1) + 0.2(1.0) \end{array} \right.$$

$$V_2(PH) = -1.6$$

2. Assuming that we are acting for three time steps, what is the optimal action to take from the starting state? Justify your answer.

The convergence of the system is $PH = -3.00$, $PF = 2.00$, and $RH = 0.00$, so the optimal strategy would be to eat first as it's the highest value. The same is true if we just consider 3 iterations, however the values are different but the resulting importance is the same.