

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:

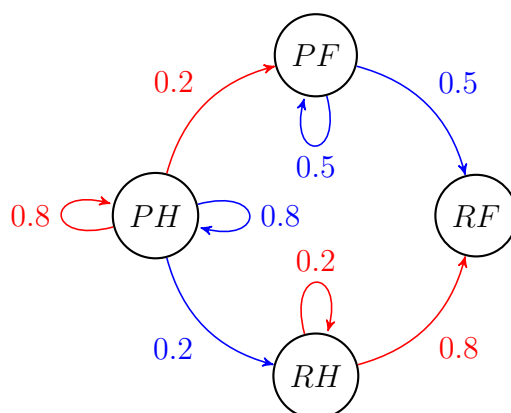
State Space {PoorHungry, PoorFull, RichHungry, RichFull}
 Actions {Eat, Play}
 Initial State PoorHungry
 Terminal State RichFull

| 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 | -0.6 | -1.2 | -1.8 |
| PoorFull | 0 | 3.0 | 2.4 | 1.8 |
| RichHungry | 0 | 4.0 | 3.4 | 2.8 |
| RichFull | 0 | 0 | 0 | 0 |

Color denotes action taken, as defined in the above chart.

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 optimal solution would be to **Play** until reaching the state *RichHungry* as it has the higher expected value of return. It converges to a value of 2.8 while **Eat** action and heading to *PoorFull* has an expected value of 1.8.