

CS188 Spring 2014 Section 4: MDPs

1 MDPs: Micro-Blackjack

In micro-blackjack, you repeatedly draw a card (with replacement) that is equally likely to be a 2, 3, or 4. You can either Draw or Stop if the total score of the cards you have drawn is less than 6. Otherwise, you must Stop. When you Stop, your utility is equal to your total score (up to 5), or zero if you get a total of 6 or higher. When you Draw, you receive no utility. There is no discount ($\gamma = 1$).

1. What are the states and the actions for this MDP?

score(state) < 6: A = {Draw, Stop} S = {0,2,3,4,5}
otherwise: A = {Stop}

2. What is the transition function and the reward function for this MDP?

s	a	s'	R(s,a,s')	P(s',a,s')
s	Stop	same state	score(s)	1
0	Draw	2,3,4	0	1/3, 1/3, 1/3
2	Draw	4,5,0	0	1/3, 1/3, 1/3
3	Draw	5, 0	0	1/3, 2/3
4	Draw	0	0	1
5	Draw	0	0	1
End-game	None	None	None	None

3. Give the optimal policy for this MDP.

state transition rewards?
or state accumulation
rewards?

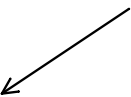
s	pi(s)	V ^{pi} (s)
0	Draw	1/3*(3+3+4)=3.3333
2	Draw	0+1/3*4+1/3*5 = 3
3	Stop	3
4	Stop	4
5	Stop	5

4. What is the smallest number of rounds (k) of value iteration for which this MDP will have its exact values (if value iteration will never converge exactly, state so).

3 ITERATIONS

- Value iteration:
- Initialize with 0 values
 - Do until convergence
 - For all states
 - Consider $\text{argmax}(w.r.t\ a, V(s))$
 - using old values

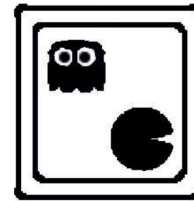
Policy has already converged!



s	V ₀	a	V ₁	a	V ₂	a	V ₃
0	0	-	0	Draw	3	Draw	10/3
2	0	Stop	2	Draw	3	Draw	3
3	0	Stop	3	Stop	3	Stop	3
4	0	Stop	4	Stop	4	Stop	4
5	0	Stop	5	Stop	5	Stop	5

2 Pursuit Evasion multi-agent sequential decision process

Pacman is trapped in the following 2 by 2 maze with a hungry ghost (the horror)! When it is his turn to move, Pacman must move one step horizontally or vertically to a neighboring square. When it is the ghost's turn, he must also move one step horizontally or vertically. The ghost and Pacman alternate moves. After every move (by either the ghost or Pacman) if Pacman and the ghost occupy the same square, Pacman is eaten and receives utility -100. Otherwise, he receives a utility of 1. The ghost attempts to minimize the utility that Pacman receives. **Assume the ghost makes the first move.**



For example, with a discount factor of $\gamma = 1.0$, if the ghost moves down, then Pacman moves left, Pacman earns a reward of 1 after the ghost's move and -100 after his move for a total utility of -99.

Note that this game is not guaranteed to terminate.

1. Assume a discount factor $\gamma = 0.5$, where the discount factor is applied once every time either Pacman or the ghost moves. What is the minimax value of the truncated game after 2 ghost moves and 2 Pacman moves? (Hint: you should not need to build the minimax tree)

Pacman never loses, Utility = Summation(from $i=0$, to 3, of $1 \cdot \gamma^i$) = $1 + 0.5 + 0.25 + 0.125 = 1.875$

2. Assume a discount factor $\gamma = 0.5$. What is the minimax value of the complete (infinite) game? (Hint: you should not need to build the minimax tree)

Pacman never loses, Utility = Summation(from $i=0$, to infinity, of $1 \cdot \gamma^i$) = 2

3. Why is value iteration superior to minimax for solving this game?

minimax has infinite depth. never converges
value iteration will converge probably pretty quickly

4. This game is similar to an MDP because rewards are earned at every timestep. However, it is also an adversarial game involving decisions by two agents.

Let s be the state (e.g. the position of Pacman and the ghost), and let $A_P(s)$ be the space of actions available to Pacman in state s (and similarly let $A_G(s)$ be the space of actions available to the ghost). Let $N(s, a) = s'$ denote the successor function (given a starting state s , this function returns the state s' which results after taking action a). Finally, let $R(s)$ denote the utility received after moving to state s .

Write down an expression for $P^*(s)$, the value of the game to Pacman as a function of the current state s (analogous to the Bellman equations). Use a discount factor of $\gamma = 1.0$. Hint: your answer should include $P^*(s)$ on the right hand side.

$$P^*(s) = \operatorname{argmax}_{a \in A_P(s)} (R(N(s, a)) + \operatorname{argmin}_{a' \in A_G(N(s, a))} (N(N(s, a), a') + P^*(N(N(s, a), a'))))$$