

Probabilistic Foundations of Artificial Intelligence

Problem Set 6

Nov 28, 2014

1. Hero in the maze

Consider the following problem related to probabilistic planning. You are a hero H who is being chased by a ghost G in a maze.

- (i) Suppose the maze is a simple (infinite) chain of nodes, each node labeled with a number (from $-\infty$ to ∞): H starts at 0, G starts at -2. H always tries to move away from G, but only succeeds with probability p , and with probability $1 - p$ gets stuck (i.e., with probability p , H moves 1 step to the right, from node i to $i + 1$, and gets 1 unit of reward; with probability $1 - p$, H doesn't change its location and gets 0 units of reward). G always chases after H and never gets stuck. If G catches H, H incurs -10 reward in the timestep in which it got caught (and 0 reward in all subsequent time steps). Both G and H move simultaneously. Write down the state space with the transition probabilities. For a discount factor γ , what is the expected long term future reward as a function p and γ ? Calculate its value for $p = .9$ and $\gamma = .95$. *Hint: You may want to consider the relative positions of H and G instead of their absolute positions when choosing your state representation.*
- (ii) Now, suppose the maze is a "T", i.e. an (infinitely large) tree, where only one node, the starting node of H, has degree 3, all other nodes have degree 2. In the first round, H has the choice of either moving "right" and being chased (the same as above); or moving "down" and not being chased. If H moves "down", it will also get stuck with probability $1 - p$ like above, but only incur reward 1/2 for each step moved (which happens with probability p). In all subsequent actions, H continues to (attempt to) move in the same direction as in the first round (i.e., once it decides to move right, it has to continue to move right etc.) What is the expected long term future reward in this case, as a function of p and γ ? Calculate its value for $p = .9$ and $\gamma = .95$. For these values of p and γ , which initial action should H take? For a value of $\gamma = .95$, give an explicit rule on how H should choose its initial action as a function of p . Compute the critical (decision-relevant) values of p (you may have to do this numerically).

2. Policy iteration

Consider an undiscounted MDP having three states, (1, 2, 3), with rewards -1, -2, 0, respectively. State 3 is a terminal state. In states 1 and 2 there are two possible actions: a and b . The transition model is as follows:

- In state 1, action a moves the agent to state 2 with probability 0.8 and makes the agent stay put with probability 0.2.
- In state 2, action a moves the agent to state 1 with probability 0.8 and makes the agent stay put with probability 0.2.
- In either state 1 or state 2, action b moves the agent to state 3 with probability 0.1 and makes the agent stay put with probability 0.9.

Answer the following questions:

- (i) Draw the MDP described above. What can be determined *qualitatively* about the optimal policy in states 1 and 2?
- (ii) Apply policy iteration, showing each step in full, to determine the optimal policy and the values of states 1 and 2. Assume that the initial policy has action b in both states.
- (iii) What happens to policy iteration if the initial policy has action a in both states? Does discounting help? Does the optimal policy depend on the discount factor?