

CS558 - INTRO TO

ARTIFICIAL INTELLIGENCE

NAME : MADHAN THANGAVAL

B ID : 200814916

POB : mthangaval

HOMEWORK - 3

② Q-LEARNING

Given:

$$\text{Learning rate} = 0.1 \quad [\times]$$

$$\text{Discount factor} = 0.9 \quad [\gamma]$$

According to α learning, we know that,

$$Q(s,a)_{\text{new}} \leftarrow (1-\alpha)Q(s,a)_{\text{old}} +$$

$$\alpha [R(s,a,s') + \gamma \max Q(s',a')]$$

We need to find the Q -value for each of the actions as mentioned in the question.

$$+ 0.1(0.9 \cdot 0) \rightarrow 0.1 \cdot 0 = 0$$

$$+ 0.1(0.9 \cdot 0) \rightarrow 0.1 \cdot 0 = 0$$

For Action -1 :

$$(s_2, north, s_3, r = -0.1)$$

$$\text{Here } \gamma = 0.9, \lambda = 0.9, R = -0.1$$

∴ using the above formulae

$$Q_{new}(s, a) \leftarrow (1 - 0.1) * 0 + 0.1 [-0.1 + 0.9(0.5)]$$

$$Q_{new}(s, a) \leftarrow 0.1 [0.35]$$

$$Q_{new}(s, a) \leftarrow 0.035$$

For Action -2 :

$$(s_3, east, s_4, r = 0.1)$$

$$\gamma = 0.9, \lambda = 0.9, R = -0.1$$

∴ using the above formulae,

$$Q_{new}(s, a) \leftarrow (1 - 0.1) * 0 + 0.1 [-0.1 + 0.9 + 0.9]$$

$$Q_{new}(s, a) \leftarrow 0.071, \quad \text{initially}$$

For Action - 2:

$$(s_2, \text{north}, s_3, r=1)$$

$$\alpha = 0.1, \gamma = 0.9, R = 1$$

done(s_3) is true

∴ using the formulae

$$Q_{new}(s, a) \leftarrow (1 - 0.1) + 0.9 + (0.1)[1 + 0.9 \times 0.9]$$

$$Q_{new}(s_1, a) \leftarrow 0.81 + 0.181$$

$$Q_{new}(s, a) \leftarrow 0.81 + 0.181$$

$$Q_{new}(s, a) \leftarrow 0.991, \quad \text{final value}$$

updated Q-values are:

$$* Q(s_2, \text{north}, s_3) = 0.071$$

$$* Q(s_3, \text{east}, s_4) = 0.071$$

$$* Q(s_1, \text{north}, s_2) = 0.991$$

$$* Q(s_2, \text{north}, s_3) = 0.991$$

$$* Q(s_3, \text{east}, s_4) = 0.991$$

2. YES/ NO.

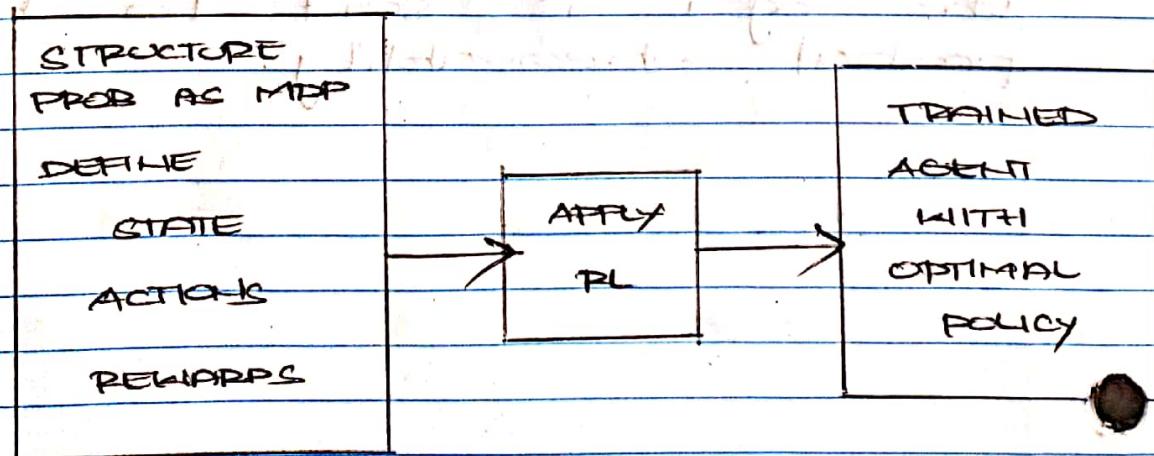
a) NO, but it assumes there exists an MDP without complete policy.

b) YES, but at sometime we don't know transition, second on both of them. [in addition to previous answer a)]

c) NO, as MDP assumes the agent directly observes the current environmental state; in this case the problem is said to have full observability. If the agent only has access to a subset of states, and the states corrupted by noise then the agent is said to have partial observability.

- d) YES, it assumes that the reward function $r(s, a')$ is available as in an MDP. (from RL definition)
- e) NO, it immediate reward is not available across the iteration and reward does not known during run time.

3.a) When we have MDP problem and RL algorithms, we will apply the RL algorithm to build an agent model and train it to find the optimal policy. Finding the optimal policy essentially solves the I problem.



like using v_1/π_1 on approx MDP and π_{true} on approx MDP. During Model-Free RL we use Q-learning and value learning to evaluate the policy.

b) When on a RL problem (and MDP algorithm ~~there~~ are given then we know the MDP then our goal is to compute v^* , π^* and π^* using value/policy iteration.

Similarly evaluating the fixed policy π using policy evaluation.

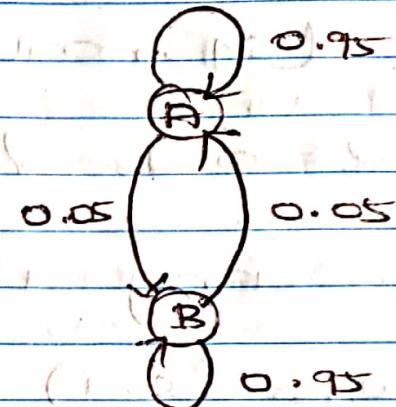
In other words, to find the optimal policy for the agent, once it has the optimal policy it simply uses that policy to pick actions from any state.

thus we will apply MDP algorithm to build an agent model and train it to find the optimal policy. Finding the optimal policy essentially solves the RL problem.

④ HMM MODEL:

a) $p(s) = ?$

s	0	$p(0 s)$
A	0	0.8
A	1	0.2
B	0	0.1
B	1	0.9



$$P = \begin{bmatrix} A & B \\ A & B \end{bmatrix} = \begin{bmatrix} 0.95 & 0.05 \\ 0.05 & 0.95 \end{bmatrix}$$

stationary distribution is the distribution we end up with P_∞ of the chain.

It satisfies

$$P_\infty(x) = P_{\infty+1}(x) = \sum p(x|x) P_\infty(x)$$

$$P_\infty(A) = p(A|0) P_\infty(A) + p(A|1) P_\infty(B)$$

$$P_\infty(B) = p(B|0) P_\infty(A) + p(B|1) P_\infty(B)$$

$$P(s) = P_\infty(A) + P_\infty(B) = 1$$

$$P_{\infty}(n) = 0.8 P_{\infty}(A) + 0.1 P_{\infty}(B)$$

$$P_{\infty}(A) \approx$$

$$2 P_{\infty}(n) = 2 P_{\infty}(B) \quad \text{--- (1)}$$

$$P_{\infty}(n) = \frac{1}{2} P_{\infty}(B) \rightarrow \text{--- (2)}$$

$$P_{\infty}(B) = 0.2 P_{\infty}(n) + 0.9 P_{\infty}(B)$$

$$P_{\infty}(B) = 2 P_{\infty}(n) \rightarrow \text{--- (2)}$$

Also,

$$P_{\infty}(A) + P_{\infty}(B) = 1.$$

$$P_{\infty}(A) + 2 P_{\infty}(n) = 1.$$

$$3 P_{\infty}(n) = 1$$

$$(1) \Rightarrow P_{\infty}(n) = \frac{1}{3}$$

$$P_{\infty}(B) = \frac{2}{3}$$

$$P_{\infty}(n) + P_{\infty}(B) = 1$$

Q.E.D.

⑥ $P(O_1=0, O_2=0, O_3=1)$

so \rightarrow stationary distribution.

$P(O_1=0, O_2=0, O_3=1)$ can be computed using the forward algorithm. The past is independent of the future, given the present.

We can compute the value by hiding each alpha (α) values using forward algo.

thus

$$\alpha_1^A = P(O_1|s) \times P(s_A) \quad (\text{B})$$
$$= 0.8 \times 0.9 = 0.72$$

$$\alpha_1^B = P(O_1|s) \times P(s_B) \quad (\text{B})$$
$$= 0.7 \times 0.1 = 0.07$$

now,

$$\alpha_2^A = 0.8 + [0.72 \times 0.95] + [0.01 \times 0.05]$$

$$= 0.8 \times [0.6845] \quad \text{Ans 1}$$

$$= 0.5476$$

$$\alpha_2^B = 0.7 [0.72 \times 0.05] + [0.1 \times [0.3 + 0.8] \times 0.95]$$

$$= 0.1 \times 0.0455 = 0.0045$$

$$\alpha_3^A = 0.2 [0.5476 \times 0.95] + [0.0045 \times 0.05]$$

$$= 0.2 \times 0.53$$

$$= 0.106$$

$$\alpha_3^B = 0.9 [0.5476 \times 0.05] + [0.0045 \times 0.95]$$

$$= 0.9 \times 0.0317$$

$$[0.0285] = 0.0285 \quad \text{Ans 2}$$

$$[0.0317] = 0.0317 \quad \text{Ans 3}$$

thus

$$P(O_1=0, O_2=0, O_3=1)$$

$$= 0[x_1^A + x_1^B] + 0[x_2^A + x_2^B]$$

$$+ 1[x_3^A + x_3^B]$$

$$= 0.306 + 0.0285$$

$$= 0.1345$$

so

$$P(O_1=0, O_2=0, O_3=1) = 0.1345$$

② Belief change after each observation of O_1 , O_2 , and O_3

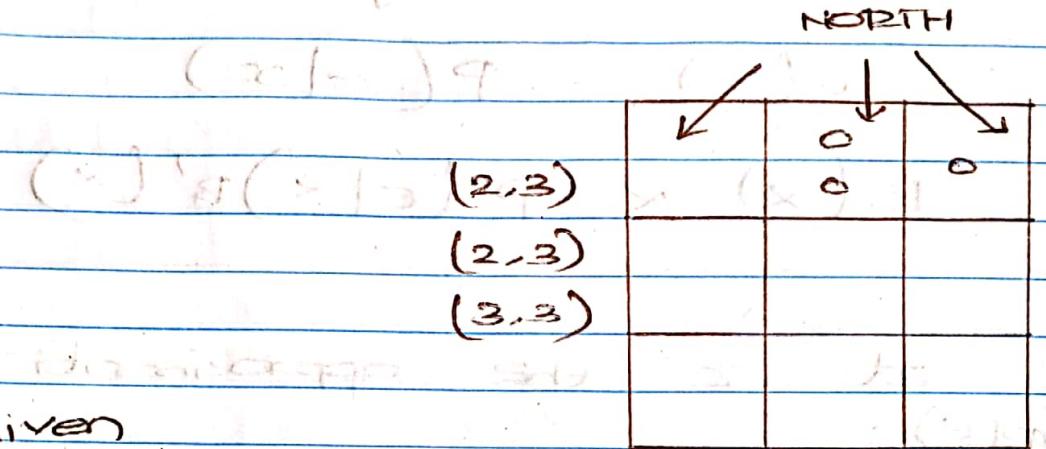
$$\text{After } O_1 = 0, = [0.72, 0.01]$$

$$\text{After } O_2 = 0, = [0.5476, 0.0446]$$

$$\text{After } O_3 = 1, = [0.306, 0.0285]$$

(5)

PARTICLE FILTER



Given

$$a) p(x) = 0.8$$

Moves \rightarrow clockwise

Here we follow the HMM structure.

- 1) Elapse (passage of time)
- 2) Weight (observation update)
- 3) Re-sample (resampled particles)

2) Elapse

$$x^1 = \text{sample}(p(x^1|x))$$

this is like prior sampling - sampling frequencies reflect the transition probabilities.

These enough samples close to exact values before and after.

2) observe (weight)

$$w(x) = p(e|x)$$

$$B(x) \propto p(e|x) B'(x)$$

It is the approximation of $p(e)$.

2) Resample particles.

In accordance with weight the particles are resampled.

It is also called likelihood weighting.

To draw the replacement we choose from the weighted sample distribution. This is equivalent to renormalizing the distribution.

ELAPSED

(2,3)

0 0

0 0

(2,2)

(2,2)

0 0

0 0

(2,2)

(3,2)

0 0

0 0

(3,2)

→ ~~and drop~~

→ ~~total height of 6
is 2 + 2 + 2 = 6~~

HEIGHT

form a
distribution
(Unbiased)

	•	•	(2,2) .2
			(3,2) .4
			(2,3) .4

EXAMPLE

New particles

		•	(3,2)
		•	(3,2)
		•	(3,2)

These are used to improve localization.

First we do the sampling, then the particle moves to the next state as it moves in the clockwise direction.

Then we observe and transition dynamics to weighted particles using probability distribution.

We weight them according to the evidence and use the weighted particles to form the distribution.

Again last we resampled the particles from weighted to unweighted particles using resampling distribution. Where the maxing weighted particles has more importance and those states are stable.

b) If we sample over and over again without any motion update, then the state (any particular) will be crowded with lot of posticks.

As the weight is relatively less than 1. Then the continuous sampler will pack all the posticks in one state.

EXAMPLE >

			...	
			...	
			...	
			...	
			...	

Here the posticks converge and forms a cluster under one state or grid of continuous weighted distribution.