

21-12-2021

Week #17 Lecture #23

Tuesday

Deep Reinforcement Learning

we won't have y , we only have $x \& z$.

first we have

$$y = f(x)$$

used x to predict y .

but now we use

x, z to predict y

Z is Reward

$$y = f(x)$$



$Z \rightarrow \text{Reward}$

x = State

In action, if it is good, we give $+\text{re}$ reward Z .

or in action, if it is not good, we give -ve reward Z .

Unsupervised Learning or unsupervised learning.

For example,

in bricks out game

then using previous Deep Neural Net techniques

we have to play a lot of games & collect

a large dataset & give it as an input & learn
steps each move right, left or stay.

This is hectic bcz it is not possible to play
all games & gather dataset.

Now:

State:

basically input for image, full screen image

in above scenario, screen, or ball position,
direction, bricks available, bricks position

Action:

paddle movement right, left, or stay

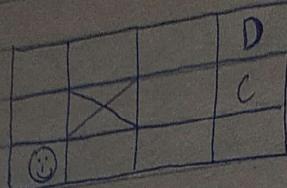
Environment:

This entire game

Reward:

usually score of game & we try to maximize it.

Example:-



D = Destination
C = Concert
X = Not reachable

Reward g could be +ve or -ve. All together 12 states.

Agent has to reach Destination.

1 episode: Agent starts from its initial position do a terminal state.

D has reward 100

C has reward -100

Rewards are also given when moving from one state to another state called Immediate Rewards

we first compute utility value of each state

which is reward accumulated till terminal state.

So utility values for each state is

88	→ 92	→ 96	→ 100
84	X	92	-100
88	84	88	84

every step reward was -4

Now agent could go up or right

dig matrix

r ₃	→	→	→	○
r ₂	↑	x	↑	○
r ₁	↑	→	↖ ↗	←

c₁ c₂ c₃ c₃

There could be multiple terminal states.

This was a Deterministic environment meaning

our ~~random~~ action taken gives a deterministic state that we intended to go.

In an undeterministic environment, we have 2 ~~ways~~ reinforcement learning.

① Model Based Learning:

Transition Model = a matrix

Starting cell	Destination	Probability
1,1	2,1	0.8 Prob
1,1	1,2	0.1
1,1	1,1	0.1

This matrix gives probabilities that to reach a destination from current cell to adjacent cells.

Compute Scores.

✓ For Deterministic environment

$$S_t = \frac{r + \sum_{\pi} (S_{t+1})}{\pi}$$

r = reward immediate next
reward state
Reward

2 loops one for all states

One loop to infinity
Until matrix values
don't change

Value Iteration Algorithm.

0	0	0	100
0	0	0	-100
0	0	0	0

↓

-4	-4	96	100
-4	-4	-104	-100
-4	-4	-4	-104

↓

-4	92	96	100
-4	-4	92	-100
-4	-4	-4	-104

but in undeterministic Environment Value Iteration becomes.

S_{t+1}

92	a	100	28
-4	b	-100	28
-4	-4	-104	28

$$\alpha S_t = 0.8 \times \underset{\text{immediate reward}}{\cancel{S_{t+1}}} + 0.1 \times a + 0.1 \times b + r$$

↓
Immediate Reward

Bellman Equation.

$$\text{Utility or Value } V_t = r_t + \gamma \sum_{s'} T(s, a, s') * v(s')$$

Gamma called Discount factor

Q values is b/w 0, 1

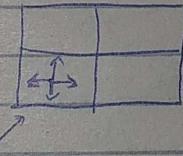
one iteration is Bellman update

Q learning:

enhanced version of Value Iteration,

here we were computing one utility value for one state

Now, for each action, we check utility value



Max value in each direction it is

ACTIONS

called Q value

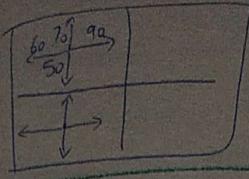
In Q learning,

we don't know
in which state we
go to from this current state

No transition Model

States	Actions			
	Up	Down	Left	Right
s ₁				
s ₂				
s ₃				
s ₄				
s ₅				
s ₆				

Model Free Reinforcement learning



Now, we have

$$Q(s, a) = Q(s, a) + \alpha \left(r + \max(Q(s', a)) - Q(s, a) \right)$$

α = learning Rate, start from a larger value & then reduce it as learning has been completed.

s = current state

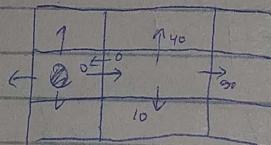
a = action

$$\max(Q(s', a))$$

↓ q₀ will be returned so Q value of
 next state is taken
 & we add immediate
 reward & we subtract it
 from current state's Q value

60 70 90
 ↓
 50

Now, Example.



$$Q(s_{21}, \text{right}) = 0 + \alpha [-4 + \max(4, 50, 0) - 0]$$

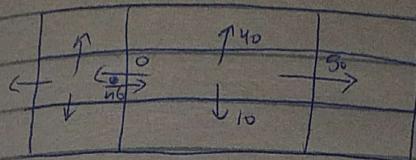
$$= 0 + \alpha [-4 + 50 - 0]$$

$$= 0 + \alpha [46 - 0]$$

$$= 0 + \alpha (46)$$

$$Q(s_{21}, \text{right}) = \alpha * 46$$

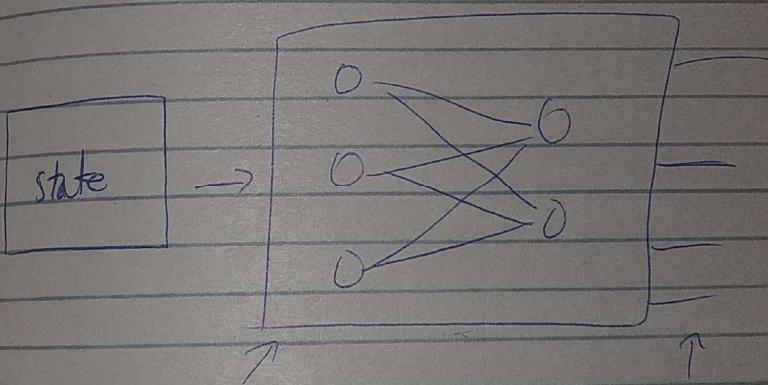
Now, we have



Problems:-

As No. of States or No. of Actions increase, this algorithm fails. So we go back to Neural Net.

Solution:



Act as Q table

Predicts Every Action Q value

we learn our Q table

Then whatever is max for one action, we assign that Q value to that action.

This \hat{Q} Equation is used for loss.