

NPTEL
NPTEL ONLINE COURSES

REINFORCEMENT LEARNING

Policy Iteration

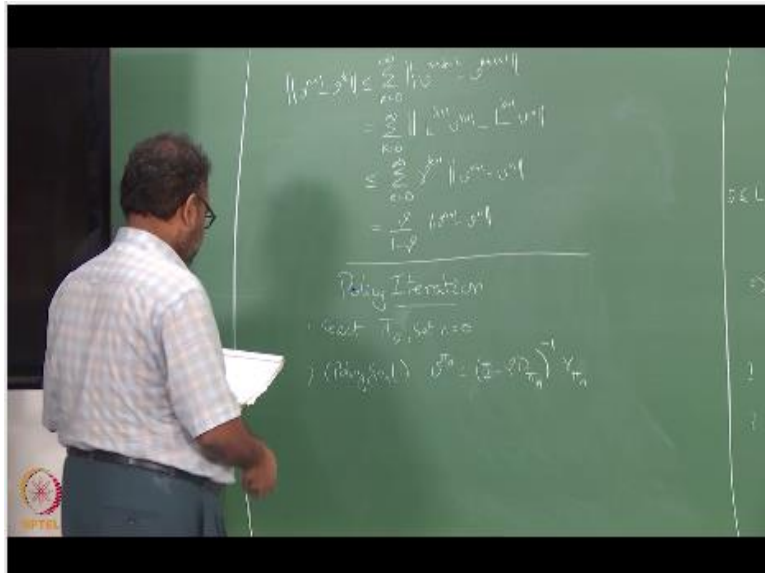
Prof. Balaraman Ravindran

Department of Computer Science and Engineering

Indian Institute of Technology Madras

So we need to do policy iteration so shall we are, should we not when we started late I know it's my mistake here oh okay, so that leaves us with 15 minutes I think, we could do it in 15 minutes, so policy iteration the idea is very simple, right so I start off with some arbitrary policy let us say π_0 may start off with some arbitrary policy π_0 I find v_{π_0} I not, right then what do, I do I be greedy with respect to v_{π_0} I will find π_1 right.

(Refer Slide Time: 02:25)

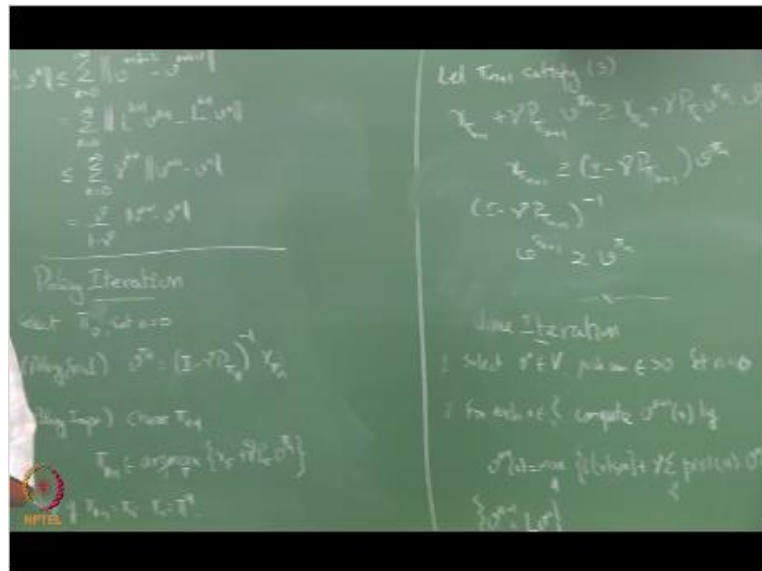


So being greedy is this right so I will take some π_0 solve find v_{π_0} , right then do this and I will find π_1 , right and then I find v_{π_1} plug that backend here right, so essentially I will plug in VN

here find π_{n+1} right, so keep going until I converge and convergence here is more straight forward if you stay with the same policy you converge where unusually convergence is much faster the rate of convergence much faster with policy, iteration but, for variety of reasons value iteration remains a very favorite algorithm for many war people don't like policy iterations for, some reasons.

So but people get the idea for the policy iteration is right you can go implement policy patients rather easy, right so first so let's of arbitrary sorry I hate going back and forth select an arbitrary π_0 , right second step is called, policy evaluation so where you say, refined v^{π_0} right so policy evaluation is just find v^{π} right so we can just do this by solving the system of equations I can just take the inverter inversion in the inverse of that matrix and solve it right or if you like it you could do the iterative fashion just just like you did in value iteration instead of iterating on L I can iterate on $l\pi$, right and then I let have some stopping criterion for stopping it but I can do this as well next thing is called policy.

(Refer Slide Time: 02:25)



Policy improvement okay this is a second step is called policy evaluation third step is called policy improvement, just like we mentioned earlier it says mark Maxes then component-wise, right for every state yes you do the max, it is whatever I wrote here before erased it at the beginning at exactly that, only difference is here I am plugging in v_{π} it is not V^* right so this v_{π} and I do this argmax over π and at every step because I am checking for equality for stopping right.

So there could be multiple actions that give you the max here right, so if at step n you chose an action for giving a max like step n there is some action your chosen as policy for a state s in step $n+1$ if the same action belongs to the max set you will pick that actually you do not pick an arbitrary action from the maxim, suppose let us say going up and going right right both are optimal, right but in the previous step I had chosen going up for the state this step I should chose I should choose going up also in otherwise I will just keep oscillating between up and right up and right I might not figure out that I have stopped right.

So I shouldn't choose an action from the maximizes randomly right are arbitrarily I have to have some mechanism by which I consistently break ties so that, policies do not change when they do not have to, $\pi_{n+1} = \pi_n$ stop and declare that π_n is V^* otherwise go too, well otherwise increment n and go to step 2, right so it's rather easy, right so if the stopping condition holds can you see why π_n should be π^* , that doesn't mean anything it just means it's a fixed point, v_{π} should be π^* well it is a fixed point but it is a fixed point of L .

If you think about it $\max_{\pi} r_{\pi} + \gamma p_{\pi} v_{\pi}$ ray that operation that looks like L that is yell operating on v_{π} then if I get that π in the same beep I same π again, right so that means the value function hasn't changed, right so that means $V = v$, so that means it is the optimal value function right so if $\pi_n = \pi_{n+1}$, because of the way I am generating the sequence of π and if $\pi_n = \pi_{n+1}$ then π_n is π^* , also I haven't changed right because the apology has it changes the value function will interchange fiesta no because of the way we are generating here right remain so.

So if the value function hasn't changed and this is essentially this part of it is essentially L operating on $v\pi$ and it does not change so that essentially means, that it is the optimal value function, okay and what else do we need to show I said if $\pi_n = \pi_{n+1}$ then π_n is V^* π_n is a π^* by the way right, what else we have to show, that I have to show that happens right I will have to show that that happens so we will do that in two ways one we will have to show that so $v\pi_{n+1} \geq v\pi_n$, right.

So why is that why you think that will be greater, so basically π_{n+1} is a greedy policy according to be π_n , okay so taken $v\pi_n$ and action greedily with respect to it and I get π_{n+1} , so that is what I mean like π_{n+1} satisfy 3 then, right this is clear right so because π_{n+1} satisfies 3 so that is being greedy with respect to v_n well π_n is some arbitrary policy which is the same as this then that will give the same answer otherwise this will be greater, because this is the one obtained by maximizing this expression, wait I took this expression I am maximizing it and whatever gives me the maximum I am saying is π_{n+1} , right.

So this is the maximum on the left hand side so the maximum has to be greater than or equal to whatever value I will get by plugging in some other π_n here because this is the max okay this is clear right one more step, so I can say so why can I say this because this is actually equal to v_n that the right hand side will be equal to $v\pi_n$ because it's a fixed point of $L\pi$, so this is $L\pi$ operating the $L\pi_n$ operating on $v\pi_n$ such a fixed-point so this is $v\pi_n$ so essentially this is it and if I multiply both sides by, what do I get so that into this is what, right.

So every step I will keep improving okay so the last bit we need is to show that, there are only a finite number of policies that you can search through, if you are having a finite MDP finite number of deterministic policies we are searching through only deterministic policies here, only finite a number of deterministic policies you can search through if you have a finite MDP so that is basic approved it so policy iteration will converge, but every time I have to become better and

there are only finite number of things I can search through at some point so I'll have to stop, or is it okay the component ways so every component, of v is greater than sorry I didn't define it earlier I should have so this essentially means set v_{n+1} of s_1 is greater than or equal to v_1 of s_1 of s_2 s_3 for all s , right I think there are people clamoring to come in so that's why I stopped we can go out and you can ask questions.

IIT Madras Production

Funded by

Department of Higher Education

Ministry of Human Resources Development

Government of India

www.nptel.ac.in