

NPTEL

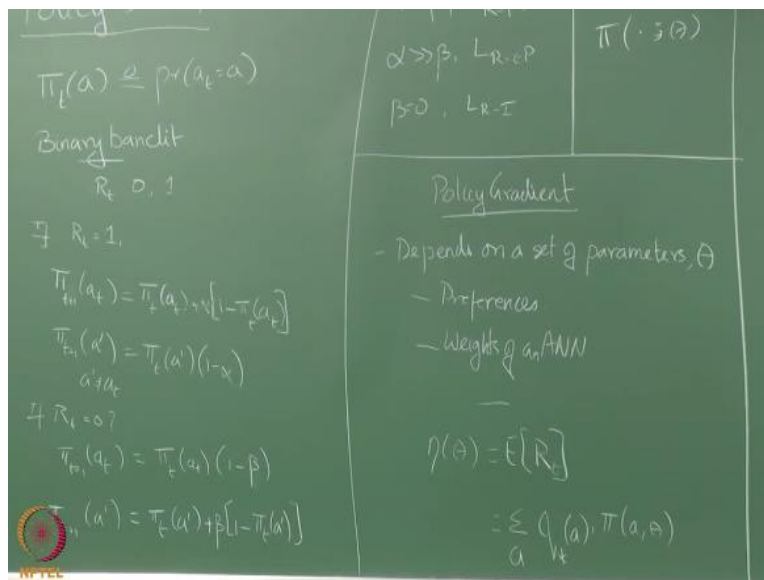
NPTEL ONLINE CERTIFICATION COURSE

REINFORCEMENT LEARNING

REINFORCE

Prof. Balaraman Ravindran
Department of Computer Science and Engineering
Indian Institute of Technology Madras

(Refer Slide Time: 00:17)



Let us denote the set of parameters me θ so I am going to come up with some kind of a performance measure for a specific choice of θ , so I will call it $E(\theta)$ right so for what is the specific choice of θ correspond to, suppose my θ is like composed of $\theta_1 \theta_2 \dots \theta_K$ so I give you values for $\theta_1 \theta_2 \theta_3 \dots \theta_K$ so what does this correspond to one, policy when I say I am looking at a specific value for θ it corresponds to one policy.

So when I say $\eta(\theta)$ that means I am evaluating that policy something okay so what is the most natural evaluation measure you can think of our policy expected pay off, right expected the expected payoff is the evaluation for that policy so yes I am going to do I am going to do that there okay, so what is the expected payoff you know what the expected payoff is, this is the expected payoff I will get for taking action A.

This is a probability that I will take action A when I am selecting actions according to θ right and so this is a standard notation I know if you are familiar with this means Π is a function that is parameterize by θ right and whatever specific value argue for θ will define what that Π is okay and then this means it is going to take some other argument okay, so this is the input to the function okay.

Which is defined by these parameters θ okay so this is how way when you put a semicolon here that means it is forever transformed okay arms are tight from taking all the notes is all okay fine yeah, so since $\Pi_A; \theta$ means that is the policy defined by the current settings of parameters θ times the $Q^* A$ right because $Q^* A$ is a two expectation for taking our may later and then some overall A I will basically get the expected reward but if I am having a deterministic policy that means essentially only one arm I will take all the other arms of probability 0 then the payoff will be Q^* of the term right okay, so now once I have this kind of a performance measure so what do I do is essentially I take θ .

(Refer Slide Time: 03:16)

$$\begin{aligned} \Theta &\leftarrow \Theta + \alpha \nabla_{\theta} J(\theta) \\ \nabla_{\theta} J(\theta) &= \sum_a q_k(a) \nabla_{\theta} v(a, \theta) \\ &= \sum_a \left(q_k(a) \frac{\nabla_{\theta} \pi(a, \theta)}{\pi(a, \theta)} \right) v(a, \theta) \\ &= E_{\pi(\cdot|\theta)} \left[q_k(a) \frac{\nabla_{\theta} \pi(a, \theta)}{\pi(a, \theta)} \right] \\ &\approx \frac{1}{N} \sum_{\eta=1}^N R_{\eta} \frac{\nabla_{\theta} \pi(a_{\eta}, \theta)}{\pi(a_{\eta}, \theta)} \end{aligned}$$

They take θ , I find the gradient of the performance okay so where the gradient is taken with respect to θ right so I look at how the performance changes with changing θ and I will try to move θ in the direction of the increase in the performance, okay. So this is called gradient ascent the people might have heard about gradient descent earlier so this is gradient descent, so I am going up the gradient, right.

If the performance improves if I change θ in some direction then I will change my θ in that direction, okay that make sense we will see this what we need more pictorial analogies right, so let us say that a one parameter θ we also have one parameter θ right and $\eta(\theta)$ let us say this nice it would be to do it right let us say that I have my performance whatever is an acceptable range of θ so my performance goes like this right.

So I initially start off with some guess for θ and that will be my initial guess for θ right then what I do, I measure this and then I compute the gradient right so the gradient is there so if I look at which direction I should move my θ , essentially I should move in this direction so I move a little bit and set that as θ_1 okay then I measure the radiant again okay still I have to move in this

direction more so I keep doing that and I think I will declare that I have reached the highest performance there.

So this is essentially how I do it okay so a couple of questions which you might ask me here, why do I have to do this in this kind of an incremental fashion why cannot I just take the derivative of that and find the highest point and settle down there, sorry I do not have the function at all so I get a good close form I do not have the function at all what is it function that is $Q^*(\Pi)$ right I know Π but I do not know Q^* .

The federal Q^* my problem is solved right I do not need to do any of these things that function itself I do not know right so how am I evaluating this function and trying to find the gradient by sampling so I am pulling a lot of farms right according to my current policy Π right and then I am trying to evaluate the policy Π and trying to find a gradient at that point right, so because everything is done through sampling.

So whatever gradient I compute at this point will not be the correct gradient right this is too simple when it has to either be that direction or this direction this is too simple right but think of a multi dimensional case when there are many parameters the right direction you will have to compute can get messed up right so instead of going in all the parameters in stuff identifying the right direction some parameters you might get the right direction some parameters you will get the wrong direction and so on so forth right.

Therefore every time you make only small steps because if we tried if we say oh I have computed the direction of the gradient I will just keep moving in the direction until my power I will take a huge step in the direction until my performance drops or something that would not work because the direction itself might be wrong so in this case I mean we have always moved in the right direction every step.

That is possible that because you are estimating the gradient you might actually can make the wrong computations and move in the wrong direction sometimes okay so what really you can hope for from these kinds of methods which are called stochastic gradient approaches okay they

are called stochastic gradient approaches because you do not know the true gradient and every time you make an estimate of what the true gradient will be and then you are using that for making moves.

I mean the most popular form of this will be called SGD which is stochastic gradient descent okay it is a very popular optimization technique nowadays but in this case you are using stochastic gradient descent but still it is the same stochastic gradient idea and so you have to be careful of how we use it and the best that you can hope to get this that in an expected sense over many updates right.

I did I did know I showed you two updates so if you make this update says this kind of updates many times in an expected sense you will move in the right direction of the gradient right so that is the kind of guarantee that typically expect for these kinds of methods okay, so is it clear now is clearer now if people are not hundred percent clear about what I am talking a stop and ask and I can elaborate at the risk of slowing down yeah.

Good point yeah so we will come to that so you can talk about this is we do not know the function itself right how are you going to estimate the gradient right, so we will come to that in a minute so we know we have a functional form for this now I am going to write down the gradient for this expression right this is expression we here I am going to write down the gradient for this expression.

So this is assuming that $Q^* A$ because it is a parameter of the problem right it is nothing to do with the policy made as far as your policy parameterization is concerned Q^* is a constant right because it is a it is something to do with the problem so I can I do not have to do this I just have to differentiate this with respect to θ right, so now I am going to do some hand waving. So what did you do here a multiplied and divided by Π , okay.

So for this to work what do I need Π is non zero for all a so this is one condition that you need for all of this any kind of Π that you assume okay has to be nonzero for all A I just assume some nonzero problem, now if you think about this particular expression right I will put a bracket

around that right so what does this remained you off, as it look like an expectation computation I am looking a the expected value of some function right.

Were I am taking samples according to probability Π right this looks like an expected value rate and summing over all possible outcomes of A all possible values here right that is the probability so I am essentially taking the expectation taken according to Π of is it fine right. So everybody on board with it right so that I can write their expectation oh okay sorry about the font, so the expectation is taken with respect to Π , okay.

Is it clear so now we know how to make an estimate of the expectations right how do you do that we draw a sample and then take an average right so we can do that so essentially what I am going to do is pull an arm that is a sample right when you say draw a sample according to Π but essentially means I have to pull an arm and what happens when I pull an arm, I get a reward right but the reward does not figure in here.

The expectation of the reward directly figures in here a so I can actually write this as a double expectation you know mean the expectation is taken over not only Π but also over the process that generates the reward Q^* right this also the process that generates the reward so my sample here is going to consist of RT right is it clear I mean this is a very subtle thing a ripple or with me here on this is it clear.

So Q^* itself is already an expectation right so instead of because I do not know this expectation I will have to estimate this expectation as well and that can be done by just taking the RT 's so what I will do now is at every time t I will sample or now I will pull the arm figure out what the reward is right but is that sufficient, no I have to do this right I just cannot use the RT directly I have to use this also right.

So I will take the gradient of Π so valuated at A divided by the so this is my one sample right this is because the expectation am taking is of this quantity so the sample that I draw is actually this right and I have to do this over I do not know some number n Π 80 you are right, defense oh no I

am not using t I am sorry if I promised you that I will use n for discrete time and T for convenience time.

So we are talking about discrete times here it has to be n I do not know if i use T here sorry about that but all of you a raise in your notebook that T and write down N , for this part of see things sorry about that right so this is clear is that anything we see right, so I need to have $1/n$ that so this gives me the radiant right so this is fine oh it is still okay yet not because this is actually a estimation okay.

So I cannot put equal to that it is only an approximation rate so can we compute this, we can right because we know Π we are the ones who determine what Π should be right because at the beginning I choose a parameterize representation that I can choose whatever suppose it is off max well I know how to take the derivative of soft max at or if it is a continuous actions I can just use a say a Gaussian right I know how to take the derivative of a Gaussian with respect to its parameters right.

Or if it is a multinomial what do I do you know how to take the derivative of a multinomial as well right so the derivative is very simple so I know what form of the function I choose it so I can go ahead and take the derivative, there are a couple of reasons for it so estimating q store need not always be the right thing to do okay because I have not introduced all the complications and the hesitating whether I should go there or not.

But the thing is so the Q^* function sometimes turns out to be incredibly complex okay as a function of the state space okay not necessarily here way as a function of the state space q turns out to be incredibly complex where as a direct representation of the policy turns out to be a lot more simple right for example there are many inventory control problems where the policy essentially turns out to be if the inventory is below a certain level you buy inventories above a certain level you do not buy okay.

But then the value function itself becomes a very complicated function of how many of the items are left of different types and so on so forth so the function itself becomes a lot harder to learn so

there are instances like this where this is easier to learn right and there are a couple of other cases especially when you are talking about very large problems for example we will see that policy gradient approaches generalize very easily to problems with continuous actions.

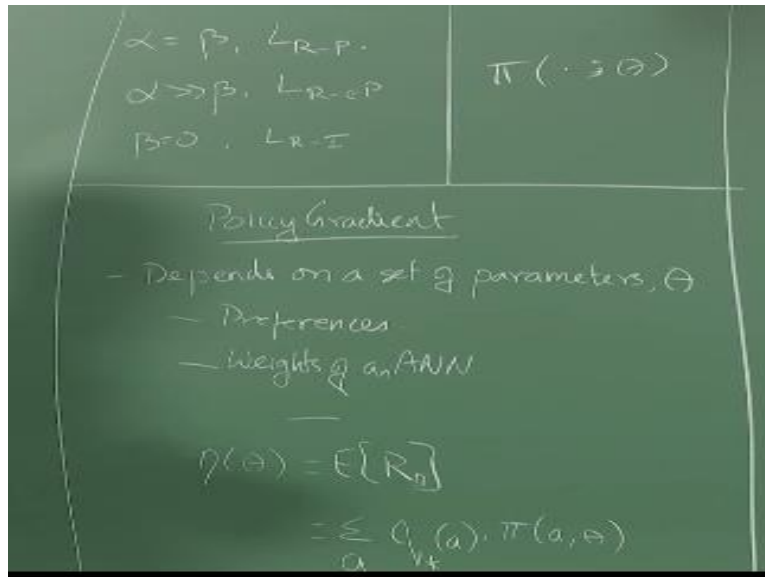
So whether individual action super if you are talking about value function methods right where I typically store a value for each action now I am going to have continuous actions it is not clear how to handle it there are ways of handling it but it turns out that they are not that well behaved as policy gradient approaches are and the action space is continuous right so there are many reasons why you want to consider the contractions I mean policy critique metrics.

Primarily I would say the Prime one of the primary motivations is to handle continues actions, okay. So this is some kind of a batch updates, what because the expectation is taken with respect to the Π now you could do more simplifications of that if you want to write but what is your question again, now but I need to take this expectation with respect to the sample the probability distribution Π right.

So and I cannot separate this I cannot write this out I saw some really cannot simplify this further because my Q^* is here right so this is the expectation of this whole quantity so I really cannot split this apart and say that okay I will just simplify this expectation and take only the expected value of Q^*A can I, I cannot take $E(AB)$ is not $A(EFB)$ is it in, A is only a constant rate not otherwise I am not sure that we can simplify this further right.

So what we are doing here is some kind of a batch mode so if you think of how we are doing this sampling so I am taking some θ fixing it pulling the arms multiple times right fixing θ pulling arms capital n times and then taking this average right and then or computing the gradient and then using that to change the parameters and this is one way of doing it another way of doing it to stay that okay I will have a θ I will pull an arm I will compute the gradient. And I will change the parameters wait so you can have a some kind of an incremental version.

(Refer Slide Time: 21:59)



So where I am saying at every step I will change my parameters by some fraction delta T I will do this at every step right sorry and so this is what I am doing so what I am essentially will end up doing is θ will be θ plus delta θ , θ_{n+1} will be θ_n plus delta θ_n right so this basically and I will do you said every time you for every time I pull an arm I will change this right I mean that is essentially the gradient of the lawn it turns out this simplifies your life a lot this observation simplifies your life a lot when you are trying to actually solve the solve problems later right now you can immediately think of soft max becoming very easy to differentiate we have now taken log and there is also lots of E's there right.

So all the e power will thing will vanish so it will become lot easier to differentiate for you right so what I am going to do now is to include another term here rather than arbitrarily called BN which we refer to as the reinforcement baseline right and adding the BN should not affect the whole process that is what the BN should not really be a function of A yeah θ basically so it should not be a function of the action I take right.

If it is a function of the action I take then it becomes dependent on θ so as long as it is not a function of the action I take I can keep adding a I can add this thing it is called the reinforcement

baseline right so one way of thinking about it is I am when I get a reward right I do not know if the reward is a good reward or a bad reward so giving a baseline and saying that if you are above this baseline it is a good reward if you are below this baseline it is a bad reward.

It kind of allows you to calibrate the rewards right one typical way in which this reinforcement base line is used is to just take the average of all the rewards I have obtained so far later this is all the reverse regardless of which are my play I keep accumulating the rewards right now when I play a particular arm if it is above this average then it is a good down which below this average then it is a bad arm right.

So valid after yeah I eventually it will eventually it will be all fine right so initially you might be making some mistakes because the initial fury will B wrong right so you might actually think a good arm is bad and a bad arm is good I mean depending but eventually it will all work out does not matter so if you think about if the reward is good right then I will go in the direction of the gradient right but if the reward is bad and I will go in the opposite direction of the gradient right.

So in some sense that is a little waving hear that but then it is fine it turns out that adding the reinforcement baseline makes the convergence behavior of this algorithm a little more stable but it is not necessary you can if you are going to implement this you can do it without the reinforcement baseline also right so this is called the and this term is sometimes called the call the characteristic eligibility.

So what is the characteristic eligibility it essentially tells you, if you do that then you get into a critic I will talk about that later and talk about that later yeah sure yeah you could do that right and that is essentially the one of the motivations for getting attacked or critical rhythms, so that SN requires you to maintain Q star estimates as well as the Π estimates as well as a Π representation.

So the Π would be that critic and the Q star would be the actor other way around the Π would be the actor the Q star will be the critic so that is another class of algorithms right which address some of the drawbacks of policy gradient approaches and so I will talk about it but somewhere

halfway down the course not now so after this I am going to go back and start talking about the full RL problem and at some point I will come back and do a TRO critic yeah.

So his question was why do I use RN right why cannot I plug-in q there right you could and yeah and typically when you do that it performs better the algorithm it is actually a it reduces variants significantly at the cost of adding significant bias and there are other issues you know you well you know or other issue is going to actor critic but I will talk about it later right and we are going back.

So why is that called the characteristic eligibility so if you think about it so it essentially chooses which of the θ are most responsible for a change at that point right so that is the term that determines that is the term that is the θ dependent term in your updates right and it tells you which θ is more responsible for change that is happening here there are three or four different that I have right.

So whichever has the highest gradient whichever direction you have the highest gradient that is a direction where the largest change is going to happen right, so it is essentially this telling you which θ is more eligible to receive the update that is why it is called the characteristic eligibility I did that make sense to people I think i have lost more people today than in the previous classes no great.

So this whole thing right this way of doing this incremental version of doing this update is actually called the called the reinforce algorithm the first proposed by this guy called Williams in 88, Williams suppose reinforce in think in 1988 and in fact he proposed this in the context of neural networks right so he when he came up with this algorithm neural networks were still ruling right.

So he proposed this in the context of neural networks he said that your θ is our weights of a neural network right and then he came up with this update and then the biggest contribution he did was he said that a this looks like a really crazy update and I mean how is this going to work and initially people are very skeptical he actually showed that in an expected sense that even

though I am doing this update after just pulling one arm in an expected since the gradient will be in the right direction right.

If I repeat this experiment multiple times and I watch how the weights evolve right they will actually move in the same direction as it would have happened if I had computed the correct gradient and then taken steps in that direction so essentially he established that for this reinforce update and it since then it has kind of been the one sole hope of convergence for all kinds of complex function approximations and RL okay.

So it works really well but it is extreme extremely, reinforce is extremely slow because the variance is very high so when I want to move on from this I will talk about the other problems of reinforce here is a horrible part about the paper so reinforce is actually an acronym there is an expansion each letter stands actually for a word he came up with a name for the algorithm which actually shortened to reinforce.

And then that started the trend for convoluted names in him since the RL community okay anyone knows what the expansion of real food you know what the expansion of reinforces forgot okay, I cannot for the life of me I refuse to memorize what reinforce transform, great so let us do some special cases of reinforce right I want to consider a binary bandit.

(Refer Slide Time: 33:26)

$$\begin{array}{l} \alpha = \beta, LR = P. \\ \alpha \gg \beta, LR = CP \\ \beta = 0, LR = I \end{array} \quad \pi(\cdot; \theta)$$

Two Actions, arb. rewards

$$\pi(a, \theta) = \begin{cases} \theta & \text{if } a = 1 \\ 1 - \theta & \text{if } a = 0 \end{cases}$$
$$\frac{\partial \ln \pi}{\partial \theta} = \frac{a - \theta}{\theta(1 - \theta)}; \quad \alpha = \beta = \theta(1 - \theta), \quad b = 0$$
$$\Delta \theta_n = \beta(a - \theta) \cdot R_n$$

Or let us say I take it banded with two actions okay but they can have arbitrary rewards here is the parameterization I choose this is essentially the Bernoulli thing right I have two actions so some probability θ I will select action 1, $1 - \theta$ I will select action 0 so what do I need to do now to get that update rule we have to find $\partial \ln \Pi / \partial \theta$ so what this will be, it will be 1 VI is one it will be maybe minus 1 here.

If correct no yeah look at Log Π it is not $\partial \Pi / \partial \theta$ yeah so it will be $1 / \theta$ via is one and it will be $-1 / (1 - \theta)$ if it case zero right, can read it like that this of two day is 1 so $1 - \theta$ or minus θ will get cancelled out will get 1 by θ this up today is zero θ and θ will get cancelled out we will get minus 1 by $1 - \theta$ okay the compact way of writing it and now I am going to say that I will choose my α to be some row times $\theta \times (1 - \theta)$.

So it turns out that as long as your α is not dependent on the actual reward that you get okay it can be dependent on the θ this is what William should as long is not dependent on the actual reward you are all fine okay you will converge so that is essentially this is result so I can make it dependent on by θ it should not be dependent on the actual action you sample okay and the actual reward that you receive so as long as that is there it is fine all right.

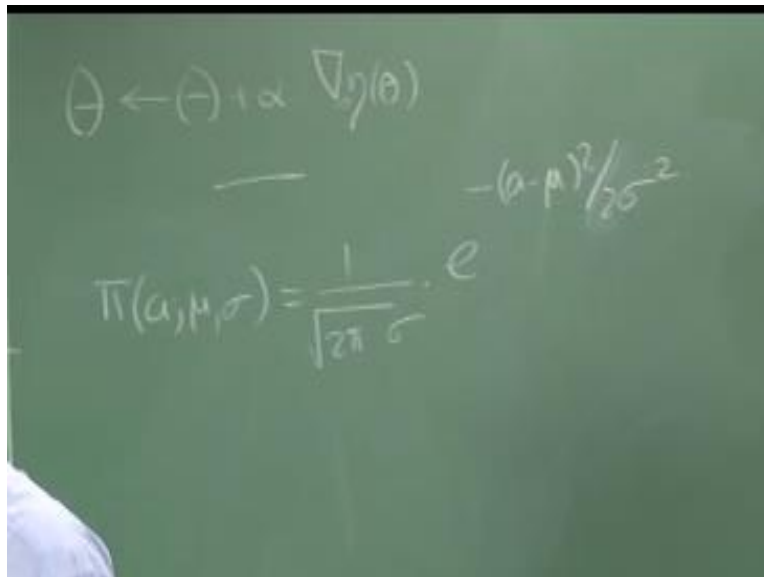
And I am going to choose my β to be 0 why am I making this specific choices so now I let me plug everything back in so what does my $\Delta \theta$ look like, so α_n which is P into $\theta_n - 1 - \theta_{n-1}$ $2 R_n - B$ this is 0 so it is $R_n - B$ into $\theta_n - 1 - \theta_{n-1}$ which will be $A - \theta / \theta \times 1 - \theta$ right so what we end up getting row okay so what is the what is the interpretation for this, let us say I take action one can and I get a reward of 1.

So what do I get P into $1 - \theta$ that is how much I change my value by so it is essentially this right $P - \theta$ right suppose I took action but what about I took action one right and my reward was one that what will I do for 0 oh there is no way there is only one parameter it is automatically my other parameter will get our district because it is $1 - \theta$ suppose I took action one and got a reward of 0 what happens no change.

Suppose I took action of 0 and got A reward of 1 what happens, minus θ into RL right so that is essentially what I have here right so this is minus θ into RL right so θ is my β if you remember right and then this is so that is essentially what I have, so this is this is what LR I right so essentially LR I is a gradient following algorithm actually it is a reinforce algorithm right so you can try to look at different choices for the α s and the different choices for Π right and try to come up with this kind of update rule right.

We will do the soft max action selection as I described in the previous class right so do that right where I replace all my Q 's with some arbitrary θ suppose I have K actions or K arms then I will have θ_1 to θ_K these are my parameters and the probability of selecting are my will be $e^{\theta_j / \beta}$ divided by summation over j $e^{\theta_j / \beta}$ right so this is my soft max expression so derive my reinforce update rules assuming the fixed pressure, right that is one I am going to give you another homework. The same thing I said one of the nice things about doing reinforce.

(Refer Slide Time: 40:14)



The image shows a chalkboard with handwritten mathematical equations. The top equation is $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$. Below it, there is a horizontal line. To the right of the line, the expression $-\frac{(a-\mu)^2}{2\sigma^2}$ is written. Below the line, the Gaussian distribution formula is written as $\pi(a; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \cdot e^{-\frac{(a-\mu)^2}{2\sigma^2}}$.

Is it allows us to handle continues actions right allows us to handle continuous actual actions need not be discreet so far way be looking at the cases or actions are discreet so I am going to use that is my probability distribution that is my policy okay give me the reinforce update rules so how many rules will I have for reinforce update rules there one for the μ one for the σ right.

Get me the reinforce update rules and yeah so try to simplify as much as possible because when you do the log μ I mean the log Π derivatives will get all kinds of we had constants coming there so you can choose your α appropriately so that some of these constants get cancelled out so that I can have again just like we did here right I chose my road to be θ into $1 - \theta$ likewise you can choose your constants in this case also write the A has to be something different so that you end up with a good nice-looking form okay great.

IIT Madras Production

Funded by

Department of Higher Education

Ministry of Human Resource Development

Government of India

www.nptel.ac.in

Copyrights Reserved