

STOCHASTIC APPROXIMATION: THEORY AND APPLICATIONS

Dr. Gagan Thope

Department of Computer Science and Engineering

Indian Institute of Science, Bangalore

Week 11

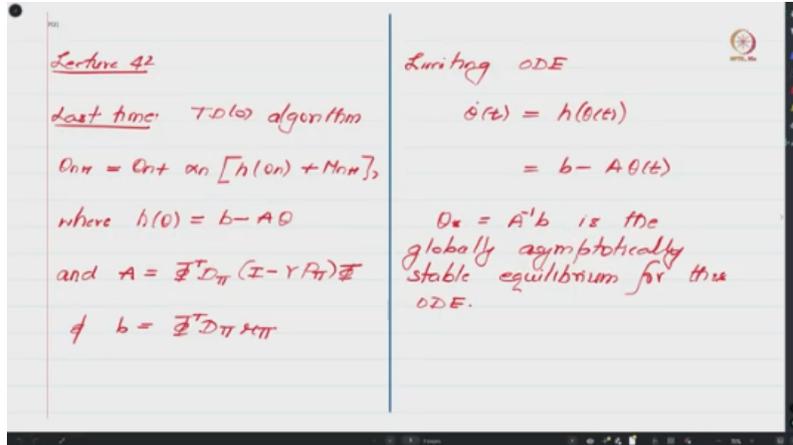
Lecture 42

Almost Sure Convergence Analysis of Temporal Difference Learning via the ODE Method

Hello and Namaste everyone. Welcome to lecture 42 of this NPTEL course on Stochastic Approximation. In this week and also in the next couple of lectures, we are looking at applications of Stochastic Approximation to reinforcement learning. In particular, we are trying to understand the asymptotic behavior of the so-called TD0 algorithm, which is used for estimating the value of a policy π in the function approximation setting. In the previous class, we looked at the limiting ODE associated with the TD0 algorithm and showed that, at least for the limiting ODE, there exists a globally asymptotically stable equilibrium, which we denoted as θ^* .

Then, our goal was to understand how good $\varphi\theta^*$ is. In particular, we wanted to compare the distance of $\varphi\theta^*$ to $V\pi$ with the distance of $\varphi\theta^{*'}$ to $V\pi$, where $\varphi\theta^{*'}$ is the projection of $V\pi$ onto the column space of φ . By this definition, one can immediately see that $\varphi\theta^{*'}$, or the projection of $V\pi$ onto φ , would have the smallest distance to $V\pi$, right? However, the TD0 algorithm—at least we have not yet guessed—would, I mean, at least we have guessed that the TD0 algorithm would not go to $\varphi\theta^{*'}$; instead, it will go to $\varphi\theta^*$.

So, right now, we are in the process of understanding how good $\varphi\theta^*$ is compared to $\varphi\theta^{*'}$. In the previous lecture, we derived some intermediate results. We will use those intermediate results to obtain the final bound that quantifies the distance of $\varphi\theta^*$ to $V\pi$ vis-à-vis the distance between $\varphi\theta^{*'}$ and $V\pi$. Let us begin the formal analysis. So, recall the TD0 algorithm has the update rule as given here, which is that $\theta_{\square+1} = \theta_{\square} + \alpha_{\square}(H(\theta_{\square}) + m_{\square+1})$, and $H(\theta) = B - A\theta$, where the matrix A has the formula given here, and B has the expression given here.



So, for this algorithm, the limiting ODE can be guessed to be $\dot{\theta}(t) = h(\theta(t))$, which, because of the nature of h of θ , will turn out to be $b - A\theta(t)$. And from this, and also the fact that your matrix A is positive definite, we managed to conclude that θ^* , equals $A^{-1}b$, is a globally asymptotically stable equilibrium with respect to the ODE that we have over here.

$$\theta_{n+1} = \theta_n + \alpha_n [h(\theta_n) + M_{n+1}]$$

$$h(\theta) = b - A\theta$$

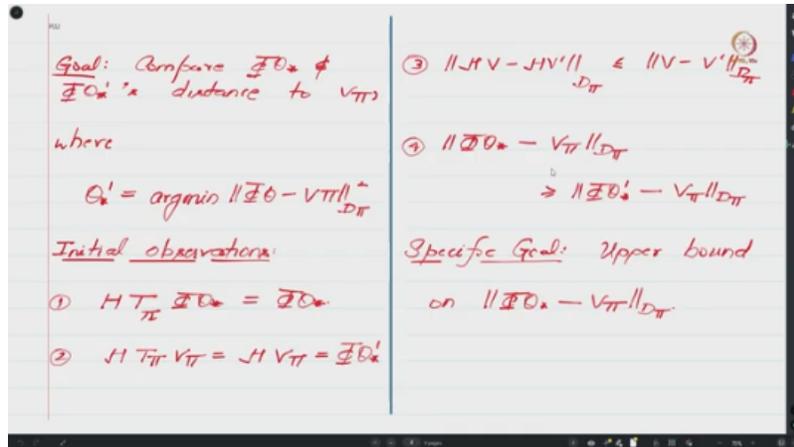
$$A = \Phi^T D_\pi (I - Y P_\pi) \Phi$$

$$b = \Phi^T D_\pi r_\pi$$

$$\dot{\theta}(t) = h(\theta(t))$$

$$= b - A\theta(t)$$

This was something that we proved in the next class—and sorry, not in the next class but the previous class—and our goal in the previous class was to somehow compare the distance of $\Phi \theta^*$ to that of $\Phi \theta^{*prime}$, and θ^{*prime} is the θ value which minimizes the distance to V_π under this norm that is induced by this matrix D_π . And we had made some initial observations in the previous class.



In particular, we had shown that this Φ_* is special in that it is the fixed point of this projected Bellman operator, that is, the operator $H T_\pi$. And furthermore, one can see that $H T_\pi V_\pi$ —one can show that V_π is actually a fixed point of the operator T_π , and hence $T_\pi V_\pi = V_\pi$. And from there, one can see that $H T_\pi V_\pi$ is actually $H V_\pi$, and $H V_\pi$ is defined to be Φ_*' . This is the way we had defined your Φ_*' . And in the last lecture, we had also seen that if you take any arbitrary vectors v and v' , look at their projections, and look at the distance between the projected quantities, then this distance is upper-bounded by the distance between the unprojected quantities, which is v and v' .

And separately, we had also seen in the previous class that $\Phi_* - V_\pi$ —this distance will be lower-bounded by the distance between Φ_*' and V_π , and this inequality holds because Φ_*' is defined to be the θ that minimizes this quantity, and hence, by this very definition, one can see that if you substitute any other vector other than Φ_*' , the distance of Φ_* to V_π would be larger than that of Φ_*' to V_π . So, our goal now is to obtain an upper bound on Φ_* — Φ_* 's distance to V_π , right? So, let us begin that computation. So, I hope you agree that, by the triangle inequality, $\Phi_* - V_\pi$

Observe that

$$\begin{aligned} & \| \Phi_{\theta^*} - V_{\pi} \|_{D_{\pi}} \\ & \leq \| \Phi_{\theta^*} - \Phi_{\theta^*}' \|_{D_{\pi}} \\ & \quad + \| \Phi_{\theta^*}' - V_{\pi} \|_{D_{\pi}} \\ & = \| H T_{\pi} \Phi_{\theta^*} - H T_{\pi} V_{\pi} \|_{D_{\pi}} \\ & \quad + \| H V_{\pi} - V_{\pi} \|_{D_{\pi}} \end{aligned}$$

Now, it is known that T_{π} is a contraction in $\| \cdot \|_{D_{\pi}}$. That is,

$$\begin{aligned} & \| T_{\pi} v - T_{\pi} v' \|_{D_{\pi}} \\ & \leq \gamma \| v - v' \|_{D_{\pi}} \end{aligned}$$

The distance to V_{π} can be upper bounded by the distance between Φ_{θ^*} and Φ_{θ^*}' and Φ_{θ^*}' 's distance to V_{π} . So, this just follows from the triangle inequality, right? And now we know that Φ_{θ^*} is the fixed point of the projected Bellman operator, and hence wherever you have Φ_{θ^*} , you can replace it by $H T_{\pi} \Phi_{\theta^*}$. And we also know that Φ_{θ^*}' is actually $H T_{\pi} \Phi_{\theta^*}$. So, this quantity is actually $H V_{\pi}$.

And, you know, from the theory of reinforcement learning, we know that V_{π} is $T_{\pi} V_{\pi}$. In other words, V_{π} is the fixed point of your Bellman operator T_{π} . Hence, this quantity equals $H T_{\pi} V_{\pi}$. So, this is exactly what we have written over here. $H T_{\pi} V_{\pi}$ in place of Φ_{θ^*}' , and wherever you have this Φ_{θ^*}' in the next expression, we will actually replace it with $H V_{\pi}$.

You will soon see why I am doing it in this fashion. And here one can see that you have H times some vector minus H times another vector. And if you recall, we had shown that the distance between the projected quantities is always upper bounded by the distance between the unprojected quantities. So by using that inequality, one can see that this expression is upper bounded by this expression over here. Is this okay? And this expression, we write it as it is. Now, you know from the theory of reinforcement learning, it is also known that your T_{π} operator is actually a contraction.

Observe that

$$\|T_{\pi} Q^* - T_{\pi} V_{\pi}\|_{D_{\pi}} \leq \|Q^* - Q^*\|_{D_{\pi}} + \|Q^* - V_{\pi}\|_{D_{\pi}} = \|Q^* - V_{\pi}\|_{D_{\pi}}$$

$$= \|\gamma T_{\pi} Q^* - \gamma T_{\pi} V_{\pi}\|_{D_{\pi}} + \|\gamma V_{\pi} - V_{\pi}\|_{D_{\pi}}$$

Now, it is known that T_{π} is a contraction in $\|\cdot\|_{D_{\pi}}$

That is,

$$\|T_{\pi} V - T_{\pi} V'\|_{D_{\pi}} \leq \gamma \|V - V'\|_{D_{\pi}}$$

So contraction in this norm. So by contraction in this norm what it means is that if you take any two vectors V and V' and look at their outputs under T_{π} . That is $T_{\pi} V$ and $T_{\pi} V'$ and look at the distance between the outputs under this metric induced by this matrix D_{π} . Then the contraction over here means that this quantity is upper bounded by γ times the distance between V and V' . So, because of the presence of γ which is the discount factor that is present in your you know MDP definition right and if you recall when I was discussing about this value of γ I had told you that this value of γ will be strictly less than 1.

Observe that

$$\|T_{\pi} Q^* - T_{\pi} V_{\pi}\|_{D_{\pi}} \leq \|Q^* - Q^*\|_{D_{\pi}} + \|Q^* - V_{\pi}\|_{D_{\pi}} = \|Q^* - V_{\pi}\|_{D_{\pi}}$$

$$= \|\gamma T_{\pi} Q^* - \gamma T_{\pi} V_{\pi}\|_{D_{\pi}} + \|\gamma V_{\pi} - V_{\pi}\|_{D_{\pi}}$$

Now, it is known that T_{π} is a contraction in $\|\cdot\|_{D_{\pi}}$

That is,

$$\|T_{\pi} V - T_{\pi} V'\|_{D_{\pi}} \leq \gamma \|V - V'\|_{D_{\pi}}$$

So, because it is strictly less than 1, what this inequality tells us is that the distance of the outputs of V and V' under this T_{π} operator will be strictly less than the distance between the inputs which is V and V' . So, this expression is true from the theory of reinforcement learning. Now, we can apply this inequality and show that this quantity

over here is upper bounded by gamma times the distance between phi theta star and V pi. right? So, this is what we end up with.

Handwritten notes on a digital notepad:

Hence,
 $\|\Phi_{\theta^*} - V_{\pi}\|_{D_{\pi}}$
 $\leq \gamma \|\Phi_{\theta^*} - V_{\pi}\|_{D_{\pi}}$
 $+ \|HV - V_{\pi}\|_{D_{\pi}}$
 Therefore,
 $\|\Phi_{\theta^*} - V_{\pi}\|_{D_{\pi}}$
 $\leq \left(\frac{1}{1-\gamma}\right) \|HV - V_{\pi}\|_{D_{\pi}}$

Thus, $\|\Phi_{\theta^*} - V_{\pi}\|_{D_{\pi}}$
 is guaranteed to be
 $\left(\frac{1}{1-\gamma}\right)$ - factor away from
 the smallest possible error.

Now, the nice thing about this quantity is that this quantity is exactly what we have over on the left hand side. Hence, we can take this quantity to the left hand side. So, you will have a 1 minus gamma and if you take the 1 minus gamma on the right hand side, we will eventually end up with an inequality of the following form which is that phi theta star's distance to V pi is less than 1 over 1 minus gamma times distance of HV to HV pi to V pi. So, I should say this is V pi here and this also is V pi.

Let me just make sure if I have V pi everywhere. Yes. So, I mean, this is also equivalent to your So, HV is precisely phi theta star prime minus V pi d pi. So, what this means is that the distance of phi theta star to V pi in this norm is at most 1 over 1 minus gamma factor away from the smallest possible error that one can incur by trying to approximate V pi within the column space of phi. So, this is the inequality that we will end up with, and this was the inequality that we were chasing, and we have finally managed to show.

Handwritten notes on a whiteboard:

Hence,

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$\leq \gamma \| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$+ \| H V_{\pi} - V_{\pi} \|_{D_{\pi}}$$

Therefore,

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$\leq \left(\frac{1}{1-\gamma} \right) \| H V_{\pi} - V_{\pi} \|_{D_{\pi}}$$

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

Thus, $\| \Phi^* - V_{\pi} \|_{D_{\pi}}$ is guaranteed to be $\left(\frac{1}{1-\gamma} \right)$ -factor away from the smallest possible error.

So, what this means is that if your gamma is, you know, let us say 0.9, right, then 1 minus gamma will be 0.1, and 1 over 1 minus gamma will be 10. So, what this means is that your phi theta star prime minus V pi, right. will be less than or equal to the distance of phi theta star minus V pi, and all of this should be under this norm. This will be less than 1 over 1 minus gamma, where you know, if you substitute gamma equals 0.9, you will end up with 10 times phi theta star minus prime minus V pi d pi, right?

So, this distance is the smallest distance that we can get, and theta star is a potential place to which your PD0 algorithm will converge. I say potential because we have not yet established that, right? But what this result is trying to tell us is that even if we manage to find phi theta star, the distance, or how bad this quantity will be, will be at most 1 over 1 minus gamma away from the best error that we can obtain. So, this is what we have managed to show.

Handwritten notes on a whiteboard:

Hence,

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$\leq \gamma \| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$+ \| H V_{\pi} - V_{\pi} \|_{D_{\pi}}$$

Therefore,

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$\leq \left(\frac{1}{1-\gamma} \right) \| H V_{\pi} - V_{\pi} \|_{D_{\pi}}$$

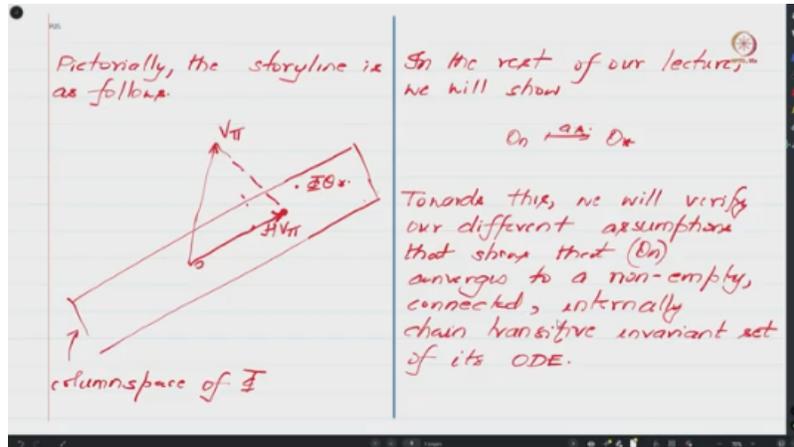
$$\| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

Thus, $\| \Phi^* - V_{\pi} \|_{D_{\pi}}$ is guaranteed to be $\left(\frac{1}{1-\gamma} \right)$ -factor away from the smallest possible error.

$$\| \Phi^* - V_{\pi} \|_{D_{\pi}} \leq \| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

$$\leq 10 \| \Phi^* - V_{\pi} \|_{D_{\pi}}$$

And pictorially, the story that we have managed to show so far can be seen as follows, right? So, here is the column space of Φ . So, this is your column space of Φ , right? And here is the origin sitting.



So, this is your origin, right? And somewhere in this column space of Φ , you will have the projection of $h v \pi$. So, $v \pi$ is the vector whose approximation you want to find in the column space of V and $H V \pi$ is that best approximation to $V \pi$ in the column space of V . Ideally, we would have liked to go over here, but the SGD algorithm that we tried working with couldn't be implemented because we don't know $V \pi$, right.

Hence, we, you know, did some approximation to that algorithm, and we ended up with the TD0 algorithm. And the TD0 algorithm, we now, you know, conjecture that it will converge to $\phi \theta^*$. And what we have shown is that this distance of $\phi \theta^*$ to $V \pi$ will be at most $1 / (1 - \gamma)$ factor of the distance between $V \pi$ and $H \pi$. So, this is the pictorial interpretation of what we have shown so far. And in the rest of this lecture, what we are going to show is that your θ_n , that is the iterates of your stochastic approximation TD0 algorithm, would almost surely converge to θ^* . And in order to show this result, we are going to make use of the different assumptions that we have studied in the first few weeks of this course.

Pictorially, the storyline is as follows

In the rest of our lectures we will show

$\theta_n \xrightarrow{A} \theta^*$

Towards this, we will verify our different assumptions that show that (θ_n) converges to a non-empty, connected, internally chain transitive invariant set of its ODE.

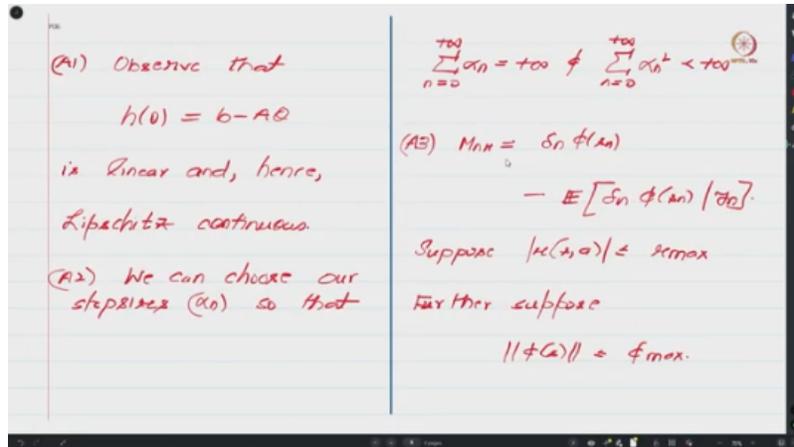
And if you recall, if we verify certain assumptions, we would be able to conclude that these stochastic approximation iterates converge to a non-empty, connected, internally chain transitive compact set of its limiting ODE. So, we are going to verify this assumption and then, you know, use this conclusion to show that θ_n indeed converges to θ^* almost surely. So, let us verify the different assumptions, and I hope you are able to remember, you know, the symbols for the different assumptions. The first assumption was A1, which required us to show that the driving function was Lipschitz continuous. In the TD0 case with linear function approximation, $H\theta$ turns out to be $B - A\theta$, which means that your H function is linear in nature.

Pictorially, the storyline is as follows

In the rest of our lectures we will show

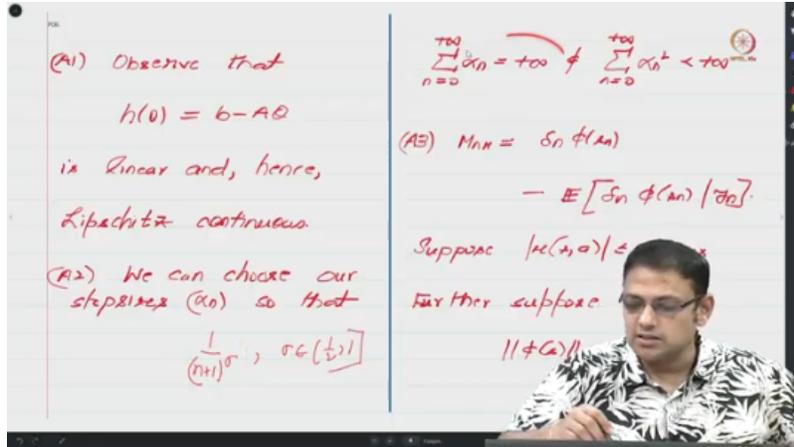
$\theta_n \xrightarrow{A} \theta^*$

Towards this, we will verify our different assumptions that show that (θ_n) converges to a non-empty, connected, internally chain transitive invariant set of its ODE.



And because it is linear in nature, it is trivially Lipschitz continuous. One can, you know, show that it is very, very easy to prove. So, hence, this assumption A1 holds true in our context. Now, the next thing that we need to verify is that your step size sequences satisfy some nice properties, which is that they should not be summable but should be square summable, right? And, you know, to ensure this condition, we will basically pick a step size choice that satisfies those conditions. And as I have told you, you know, particular examples of such step sizes include $1/n + 1$ to the power sigma, where sigma is some number, let us say between half and one, and one is close.

So, we can pick one as well. So, if you pick a step size like this, these conditions can be ensured, and because of that reason, your assumption A2 can also be trivially ensured. Now, with regards to assumption A3, we want to first show that your MN sequence is a square integrable martingale difference sequence. Furthermore, it satisfies certain linear growth conditions.



So, those conditions we will now verify. So, recall that M_{n+1} is basically $\delta_n \phi(s_n)$ minus the conditional expectation of $\delta_n \phi(s_n)$ with respect to the information that we have until time n . And what we will do is we will presume that the rewards are upper bounded by some maximum quantity that is R_{\max} . Furthermore, we are presuming that the length of your feature vectors, they are also upper bounded by ϕ_{\max} . So, let us make such assumptions like this and let us presume that R_{\max} and ϕ_{\max} are less than infinity.

$$h(\theta) = b - A\theta$$

$$\sum_{n=0}^{+\infty} \alpha_n = +\infty$$

$$\sum_{n=0}^{+\infty} \alpha_n^2 < +\infty$$

$$M_{n+1} = \delta_n \phi(s_n) - E[\delta_n \phi(s_n) | F_n]$$

So, in which case one can show that if you take the norm of $\delta_n \phi(s_n)$.

(A1) Observe that

$$h(\theta) = b - A\theta$$

is linear and, hence, Lipschitz continuous.

(A2) We can choose our stepsize (α_n) so that

$$\frac{1}{(n+1)^\gamma}, \quad \sigma_G\left(\frac{1}{n}\right)$$

(A3) $M_{n+1} = \delta_n \phi(\theta_n)$

$$= \mathbb{E}[\delta_n \phi(\theta_n) | \mathcal{G}_n]$$

Suppose $|k(\theta, \alpha)| \leq \kappa_{\max}$

Further suppose

$$\|\phi(\theta)\| \leq \phi_{\max}$$

Then,

$$\|\delta_n \phi(\theta_n)\|$$

$$\leq \|\kappa(\theta_n, \alpha_n) + \gamma \phi(\theta_n) \theta_n - \phi(\theta_n) \theta_n\|$$

$$\leq \kappa_{\max} \phi_{\max} + \gamma \phi_{\max}^2 \|\theta_n\| + \phi_{\max}^2 \|\theta_n\|$$

Hence,

$$\|M_{n+1}\| \leq 2 \kappa_{\max} \phi_{\max} + 2(1+\gamma) \phi_{\max}^2 \|\theta_n\|$$

From this, it follows that

$$= \kappa_{\max} \phi_{\max} + (1+\gamma) \phi_{\max}^2 \|\theta_n\|$$

So, recall that δ_n is this expression within the square bracket. So, if you know take the norms one can see that this quantity is upper bounded by R_{\max} in particular the absolute value is upper bounded by R_{\max} and if you multiply this thing with $\phi(\theta_n)$ this norm is upper bounded by ϕ_{\max} and hence you would end up with R_{\max} times ϕ_{\max} . And similarly this quantity and this quantity when you multiply with each other and you know make use of let us say something like the Cauchy-Schwarz inequality then one can see that this norm is upper bounded by $1 \cdot \phi_{\max}$ this norm is upper bounded by another ϕ_{\max} and hence you end up with a ϕ_{\max}^2 times the norm of this vector θ_n . And finally you know this expression similarly one can show is upper bounded by ϕ_{\max}^2 times norm of θ_n . So if you put all of them together one can see that this norm is upper bounded by $\kappa_{\max} \phi_{\max} + 1 + \gamma$ times ϕ_{\max}^2 of norm θ_n . And hence, one can conclude that.

So, if you recall M_n plus 1 is this expression plus this expression. So, if you take the norm of this, it will be the norm of this plus the norm of this thing and one can show that the norm of a conditional expectation is upper bounded by the expectation of the norm. or the conditional expectation of the norm and from that one can show that you know your norm M_n plus 1 is basically upper bounded by twice this quantity which is that 2 times R max times ϕ max plus. So, here there should be a plus 2 times 1 plus γ ϕ max square norm of θ_n . So, one can trivially show this expression. And from this, it follows that your M_n is actually square integrable.

(A1) Observe that $h(\theta) = b - A\theta$ is linear and, hence, Lipschitz continuous.

(A2) We can choose our stepsize (α_n) so that $\frac{1}{(n+1)^\gamma}$, $\gamma \in (\frac{1}{2}, 1]$

(A3) $M_{n+1} = \theta_n + \phi(A_n)$
 $= E[\theta_n + \phi(A_n) | \mathcal{F}_n]$
 Suppose $|k(x, \theta)| \leq k_{\max}$
 Further suppose $\|\phi(A_n)\| \leq \phi_{\max}$

Then, $\|\theta_n + \phi(A_n)\| \leq \| [k(A_n, \theta_n) + \gamma \phi(A_n)] \theta_n - \phi(A_n) \theta_n \|$
 $\leq k_{\max} \phi_{\max} + \gamma \phi_{\max}^2 \|\theta_n\| + \phi_{\max}^2 \|\theta_n\|$

$= k_{\max} \phi_{\max} + (1+\gamma) \phi_{\max}^2 \|\theta_n\|$

Hence, $\|M_{n+1}\| \leq 2 k_{\max} \phi_{\max} + 2(1+\gamma) \phi_{\max}^2 \|\theta_n\|$

From this, it follows that

So, you are, you know, at any given time instance, you are basically adding terms with finite norms, right? And hence, at any given time instance n , your norm θ_n is going to be bounded. And from that fact, one can show that M_n plus 1 is actually bounded. Of course, the bound can keep changing with n . But for every n , we have some bound. And

from that, one can see that because your M_n plus 1 is bounded, its expectation will also be bounded, and hence one can conclude that each M_n is square integrable.

(M_n) is square integrable
 $\mathbb{E}[\|M_n\|^2 | \mathcal{F}_n]$
 $\leq K [1 + \|\theta_n\|^2]$

$(A5) \quad h_c(x) = \frac{h(cx)}{c}$
 $= \frac{b - Acx}{c}$

Hence,
 $f_{\text{Ode}}(x) = -Ax$

Since A is positive definite, it follows that
 $V(x) = \frac{1}{2} \|x\|_A^2$

is a Lyapunov function for this ODE with respect to the origin.

Furthermore, you know, it is based on that inequality that we saw on the previous slide; one can also immediately see that if you take the square of this expression and take the conditional expectation, then one can indeed find the K such that this conditional expectation is upper bounded by K times 1 plus $\|\theta_n\|^2$. So, this tells us that your martingale noise sequence indeed satisfies the linear growth condition in terms of $\|\theta_n\|^2$. And so, this sort of verifies the first three assumptions. Now, what we need to verify is that the iterates are almost surely bounded.

And towards that, if you remember, we had looked at a sufficient condition which involves looking at the scaled ODE and showing that the origin is the globally asymptotically stable equilibrium with respect to the scaled ODE. So, towards that, we had defined this h_c function, which was h of cx by c , where c is any scalar bigger than or equal to 1 . And if you substitute, you know, this CX in the definition of H , right, one can see that we would end up with $B - ACX$, and this C will cancel off with this C . So, one can conclude that this H_c of X will basically be B over C minus AX . And if you take C to infinity, one can now see that H_∞ of X will indeed be $-AX$.

Now, we have already shown that this matrix A is positive definite. So, from this fact, one can verify that if you define a function V in this fashion—that is, V of X equals half of the norm of X squared—then one can show that. So, this function is actually a limiting ODE with respect to, you know, $\dot{x}(t) = H_\infty(x(t))$. So, sorry, I should

have been using theta here. So, let me write it in terms of theta. So, theta dot of t is equal to H infinity of theta of t. So, if you define your V of theta to be half

of the Euclidean norm squared, then one can see that this expression is the Lyapunov function with respect to this ODE with respect to the origin, which means that all solution trajectories of this ODE are guaranteed to converge to the origin. In other words, because of the existence of a Lyapunov function, one can show that the origin is indeed the globally asymptotically stable equilibrium with respect to this limiting ODE, as desired. So, this verifies your assumption A phi as well, and now, if you—so this sort of verifies all the assumptions—and hence one can conclude that indeed the iterates of your TD0 algorithm will converge to some non-empty, connected, compact, internally chain-transitive invariant set of the limiting ODE theta dot of t equals H theta of t.

$$h_c(x) = \frac{h(cx)}{c}$$

$$= \frac{b - Acx}{c}$$

$$h_\infty(x) = -Ax$$

$$V(x) = \frac{1}{2} \|x\|_2^2$$

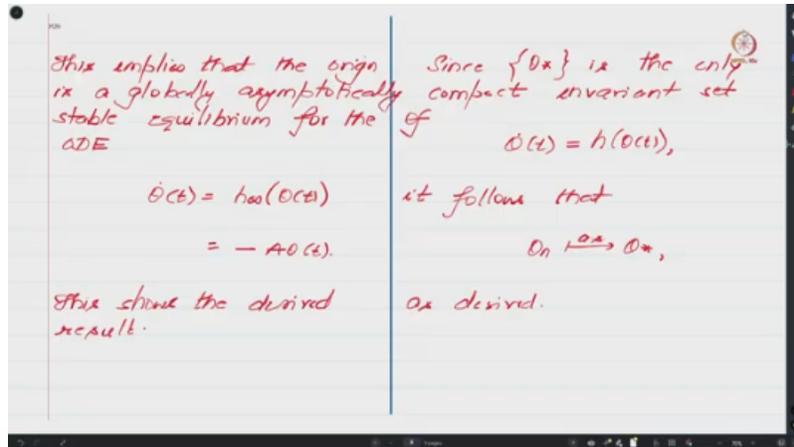
But we have already shown that the only compact invariant set of this limiting ODE is actually the singleton set containing theta star. That is the only compact invariant set. Of course, RD is also an invariant set trivially, but that RD set is not compact.

The image shows handwritten notes on a digital whiteboard, divided into two columns by a vertical line. The left column contains the following text and equations:

- (M_n) is square integrable
- $\int \mathbb{E}[\|M_n\|^2 | \mathcal{F}_n]$
- $\leq K [1 + \|B_n\|^2]$
- (As) $h_c(x) = \frac{h(cx)}{c}$
- $= \frac{b - Acx}{c}$

The right column contains the following text and equations:

- Hence, $\dot{x}(t) = h_\infty(x(t))$
- $f_{\infty}(x) = -Ax$
- Since A is positive definite, it follows that
- $V(x) = \frac{1}{2} \|x\|_2^2$
- is a Lyapunov function for this ODE with respect to the origin.
- $V(\theta) = \frac{1}{2} \|\theta\|_2^2$



So, the only other invariant set—because θ^* is globally asymptotically stable with respect to this limiting ODE, right—is basically the singleton set containing θ^* . So, from this fact and the fact that, because these four assumptions we have verified, so your θ_n has to converge to one such invariant set, and there is only one invariant set that satisfies all this criterion, one can conclude that θ_n must almost surely converge to θ^* . So, let us quickly summarize what we have done so far. We defined the TD0 algorithm, which was the algorithm that we wanted to use to come up with an approximation for V_π in the column space of ϕ , and what we have shown is that the stochastic algorithm itself converges to θ^* , and if it converges to θ^* , whatever is the limiting value that will be used as a proxy for V_π .

And separately, we have shown how good is $\phi \theta^*$ compared to $\phi \theta^*$, which is the best approximation to V_π . So, one can see that our original goal was to find a good approximation to V_π in the column space of ϕ right and we designed an algorithm to find this approximation and while we could not find the best approximation to V_π in the column space of ϕ we were able to come up with a reasonable approximation. And the nice thing about this algorithm is that it can take data that is obtained by interacting with the environment. What do I mean by that?

So, when you interact with the environment, you see that the environment is in a certain state, you take some action and then you get to see the next state and During this interaction, you may get some immediate reward, right? So, you can take these information and use this information to get an estimate of how good your strategy is,

right? So in the next week we will focus on generalization of this problem which is known as the control problem. So far we were looking at given a policy how good it is.

In particular we wanted to quantify its quality in terms of the value function. In the next class we will try to use similar principles to design an algorithm that can be used for finding the optimal strategy itself. I hope you will join me for the next week. Until then, goodbye and namaste.