

STOCHASTIC APPROXIMATION: THEORY AND APPLICATIONS

Dr. Gagan Thope

Department of Computer Science and Engineering

Indian Institute of Science, Bangalore

Week 13

Lecture 48

Q-Learning with Linear Function Approximation under ϵ -Greedy Exploration

Hello and Namaste everyone. Welcome to lecture 48 of this NPTEL course on Stochastic Approximation. This week, we have been looking at some advanced topics related to applications of Stochastic Approximation in the context of reinforcement learning. In particular, we have been focusing on understanding the behavior of Q-learning with linear function approximation and epsilon-greedy exploration. In the case of policy evaluation, we saw that the TD(0) algorithm with linear function approximation is guaranteed to converge under some reasonable conditions, and we saw that the limit point was close to

you know, the projection of V_{π} onto the function approximation space, in particular, V_{π} onto the projection of V_{π} onto the column space of your ϕ matrix. So, then we used some of the ideas from there and we developed our Q-learning algorithm for the tabular setting, and in the tabular setting as well, we saw that the Q-learning algorithm, when the behavior policy is fixed, right, converges to Q^* , which is the Q-value function of your optimal policy. And you know, things again, this convergence was guaranteed under reasonable conditions. And in the previous class, we saw or we considered the setup of Q-learning with linear function approximation, but instead of you know interacting with the environment using a fixed behavior policy, we followed the idea that is practiced in that is used in practice or applications, which is that of interacting with the environment using an epsilon-greedy policy. The idea is, whatever is your current estimate of Q^* , which is formally denoted as $\phi \theta_n$, you try to identify the greedy policy with respect to that estimate and then be epsilon-greedy according to that greedy policy, which

means that with $1 - \epsilon$ probability you act according to this greedy policy and with ϵ probability you act randomly.

And we wanted to understand the behavior of this algorithm, and in the last few slides of the previous class, I showed you some pictures which showed some very interesting behavior. That is, if you start the algorithm from some initial point, on different runs of the experiment, you converge to different places. In particular, we saw that on one of these runs, the policy or the greedy policy associated with these iterates keeps oscillating without stopping. And we saw this very interesting behavior, and we wanted to understand why that happens. In this class, we will give a formal description of why that is happening; in particular, we will try to you know help understand or come up with a technique to understand the asymptotic behavior of such algorithms.

So with that, let us begin with the formal discussion. So, as I said, we are studying Q-learning with linear function approximation and epsilon-greedy exploration. So, the idea in function approximation, in particular linear function approximation, is that we have been given this matrix Φ whose number of rows equals the cardinality of the state and action spaces and the number of columns equals d . And in this, I mean for this given Φ , the goal is to find a θ^* such that $\Phi \theta^*$ is approximately equal to Q^* , where recall Q^* is the Q value of the optimal policy. And towards that, we propose this algorithm, which is $\theta_{n+1} = \theta_n + \alpha_n (\phi(s_n, a_n) - Q_n)$, where α_n is your TD.

Lecture 4F

last time: Q-learning with linear function approximation & ϵ -greedy exploration.

Goal: Given $\Phi \in \mathbb{R}^{(S+A) \times d}$, find θ^* such that $\Phi \theta^* \approx Q^*$

Proposed Algorithm:

$$\theta_{n+1} = \theta_n + \alpha_n (\phi(s_n, a_n) - Q_n)$$

where

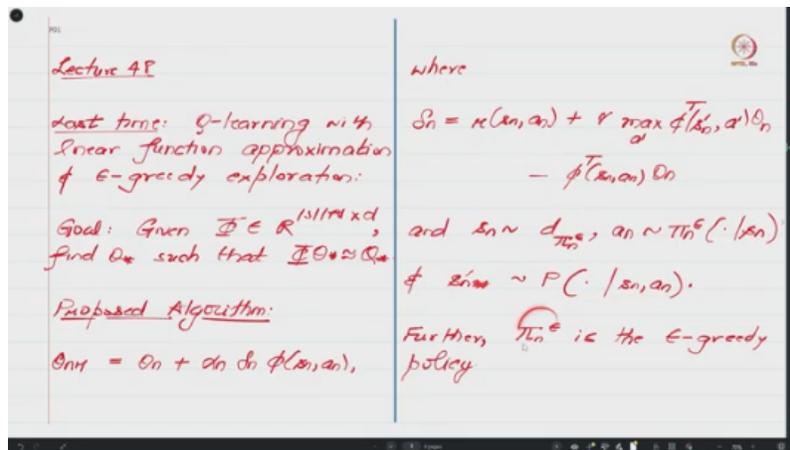
$$Q_n = r(s_n, a_n) + \gamma \max_a (\phi(s_n, a) Q_n - \phi(s_n, a) Q_n)$$

and $s_n \sim d_{T_n^\epsilon}$, $a_n \sim T_n^\epsilon(\cdot | s_n)$
 $a_{n+1} \sim P(\cdot | s_n, a_n)$.

Further, T_n^ϵ is the ϵ -greedy policy

TD error or the temporal difference error and is described as follows. It is the immediate reward plus the approximation of Q star or the current approximation of Q star at state SN prime and you take the max according to that. Minus the, you know, current approximation of Q star at state SN AN. So this is the approximation of Q star at SN prime A prime and this is the approximation of Q star at SN AN and one can think of this expression again as some approximation to Q star SN AN and this is another approximation of Q star SN AN and in this way one can see that we are taking the difference between two approximations of Q star SN AN and that is why this is known as the temporal difference error and unlike the TD0 setup

Right here Sn is sampled according to the stationary distribution that is associated with your epsilon-greedy policy and An is again chosen according to your epsilon-greedy policy and the next state again I should emphasize this is not Sn plus 1 rather it is Sn prime which is present over here. You can see that Sn prime is present. So this Sn prime is presumed to be sampled from your transition kernel. And recall that pi n epsilon is the epsilon-greedy policy. In other words, your pi n is a greedy policy.



With respect to phi theta n and pi n epsilon is 1 minus epsilon pi n plus epsilon pi n right.

$$\Phi \in R^{|\mathcal{S}| \times |\mathcal{A}| \times d}$$

$$\Phi \theta \approx Q_*$$

$$\theta_{n+1} = \theta_n + \alpha_n \delta_n \phi(s_n, a_n)$$

$$\delta_n = r(s_n, a_n) + r\phi^T(s'_n, a'_n)\theta_n - \phi^T(s_n, a_n)\theta_n$$

$$s_n \sim d_{\Pi_n}$$

$$a_n \sim \Pi_n^\epsilon(s_n)$$

$$s'_n \sim P(s_n, a_n)$$

So, this is the setup that we considered and as I told you, this algorithm had some strange behaviors. In particular, we had considered a two state two action environment right and we were using a two dimensional linear function approximation. And we saw that if you start even very close to the representation for the optimal Q-value function—in particular, we had chosen the phi matrix in this example to ensure that it is perfectly realizable, which means that the Q-star actually lies in the column space of phi. We had ensured that by keeping one of the columns of phi to be Q-star. Right. And so this is that point.

Lecture 4E

task time: Q-learning with linear function approximation & ϵ -greedy exploration.

Goal: Given $\Phi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times d}$, find θ_* such that $\Phi \theta_* = Q_*$

Proposed Algorithm:

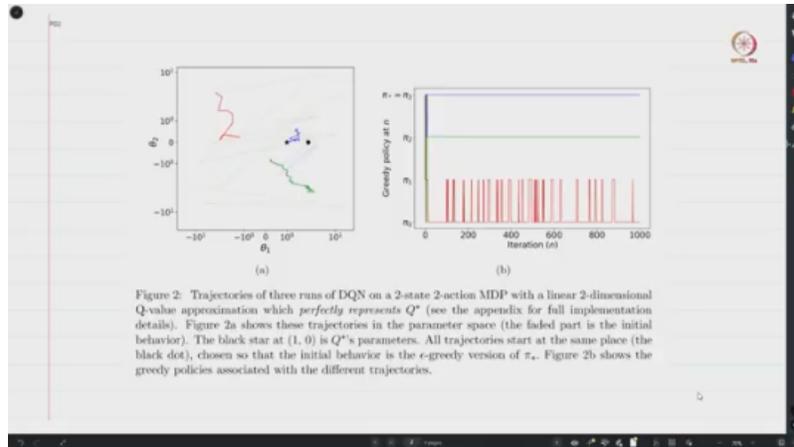
$$\theta_{n+1} = \theta_n + \alpha \delta_n \phi(s_n, a_n),$$

where

$$Q_n = r(s_n, a_n) + r \max_a \phi^T(s'_n, a') \theta_n - \phi^T(s_n, a_n) \theta_n$$

and $s'_n \sim d_{\Pi_n^*}$, $a_n \sim \Pi_n^\epsilon(\cdot | s_n)$
 $\phi^T(s'_n, a'_n) \sim P(\cdot | s_n, a_n)$.

Further, Π_n^ϵ is the ϵ -greedy policy Π_n is greedy wrt Q_n
 $\Pi_n^\epsilon = (1-\epsilon) \Pi_n + \epsilon \Pi_n^*$



So this is 1, 0. Right. And this point corresponds to the I mean, if I multiply phi with 1, 0, you will get the first column of phi. And we had said that to be Q star. And one can see that we choose a theta 0, which is very close to this.

You know, the ideal place where we would like to go, but despite starting so close to this, one can see that the noise pushes you very far away often and along different runs of the experiment you sort of reach different points. The blue trajectory ended up reaching the place where we would have liked. The green trajectory, you know, sort of settled somewhere else and the red trajectory also sort of settled somewhere else. And the interesting part of the red trajectory is that if you look at the greedy policy that is associated with the phi theta and estimates along this red trajectory, the greedy policy keeps jumping between pi 0 and pi 1 and surprisingly none of them are optimal, right? So it keeps jumping between pi 0 and pi 1.

And you know if you are in an applied setup and if you see such a behavior you may want to know you know should I choose pi 0, should I choose pi 1 and so on. But here we observe that neither pi 0 nor pi 1 are optimal but you know the greedy value associated with your phi theta actually keeps oscillating and it does not stop. You can see that the oscillation continues throughout our experiment. So the question that we had is why does this happen? Can we provide a rigorous explanation of why this is happening?

Analysis:

$$E[\delta_n \phi(x_n, a_n) | \mathcal{I}_n]$$

$$= \sum_{s, a, s'} d_{\pi_n^\epsilon}(s) \pi_n^\epsilon(a|s) P(s'|s, a)$$

$$\times \left[\kappa(x_n, a) \phi(x_n, a) + \gamma \phi(x_n, a) \max_a \phi(\bar{s}, a) \theta_n - \phi(x_n, a) \phi(\bar{s}, a) \theta_n \right]$$

$$= \Phi^T D_{\pi_n^\epsilon} \kappa + \gamma \Phi^T D_{\pi_n^\epsilon} P_{\pi_n^\epsilon} \Phi \theta_n - \Phi^T D_{\pi_n^\epsilon} \Phi \theta_n$$

where $P_{\pi_n^\epsilon}(s', a' | s, a) = P(s' | s, a)$
 if $a' = \operatorname{argmax}_b \phi(\bar{s}, b) \theta_n$
 $\neq 0$ otherwise.

So towards that let us begin the formal analysis. So as we had done previously what we will do is we will take this expression that is present in your update rule for your Q learning with linear function approximation and take its conditional expectation. And making use of the formulas for delta n and phi s n a n and also recalling the way delta like s n a n and s n prime are sampled, one can see that this conditional expectation can be represented in the following way. So, we are sampling S according to the stationary distribution associated with the Markov chain induced by this policy pi and epsilon. We pick action according to our epsilon greedy policy pi and epsilon and the next state is chosen according to this transition kernel.

So, by taking the product of these probabilities and writing an expansion of the expression that we have over here in particular the first term in delta n is R S n A n. So, if you presume that you know S n and A n equal S comma A then we would end up with R S a over here and similarly this quantity will translate to phi S a. In the same way, the next quantity will translate to gamma times phi S A times the max of this expression and the last term in your delta n expression when it is multiplied by phi will translate to something like this. And one can now try to compactly write as shown over here. In particular, one can see that This expression times these things in particular by noting that there is no S prime over here one can see that the you know expression that is over here can be compactly written in the following way that is phi transpose times D subscript pi N epsilon times R right here.

This R is a cardinality SA times one vector. This is a SA cross SA matrix and this is a D cross SA matrix. So this is how it is and your D pi n epsilon is actually a diagonal matrix whose diagonal entries are made up of, so it is a SA cross SA matrix and the SAth diagonal entry is made up of this times this expression. So this product will be present in the SAth diagonal entry or little s comma a diagonal entry of this matrix over here.

Analysis:

$$E[s_n \phi(x, a_n) | \mathcal{F}_n]$$

$$= \sum_{s, a} d_{s, a} \pi_n^a(a|x) P(s|x, a)$$

$$\times \left[x(x, a) \phi(x, a) + \gamma \phi(x, a) \max_a \phi^T(s', a') \theta_n - \phi(x, a) \phi^T(s, a) \theta_n \right]$$

where $P_{\pi_n}(s', a' | x) = P(s' | x, a)$
 if $a' = \text{argmax}_a \phi^T(s', a) \theta_n$
 $\neq 0$ otherwise

So if you take this expression and multiply it with this one can see that you can compactly write it in this way. Similarly this expression when you multiply it with this way and here you can see that there is an S prime present that leads to some expression like this where you additionally have this P pi n. So P pi n again is a S A cross S A matrix whose s prime a prime given sa right this entry will equal p of s prime given sa if this a prime over here maximizes this quantity right or it is 0, right. So if you define P pi n in the following way, one can then see that because of the presence of this max expression, right, we can write this times this in the form that is given over here.

Analysis:

$$E[\hat{\theta}_n \phi(x_n, a_n) | \mathcal{F}_n]$$

$$= \sum_{s, a, s'} d_{\pi_n^e}(s) \pi_n^e(a|s) P(s'|s, a)$$

$$\times \left[\kappa(x, a) \phi(x, a) + \gamma \phi(x, a) \max_a \phi(\bar{s}, a) \theta_n - \phi(x, a) \phi(\bar{s}, a) \theta_n \right]$$

$$= \mathbb{D}^T D_{\pi_n^e} \kappa + \gamma \mathbb{D}^T D_{\pi_n^e} P_{\pi_n^e} \mathbb{D} \theta_n - \mathbb{D}^T D_{\pi_n^e} \mathbb{D} \theta_n$$

where $P_{\pi_n^e}(s', a' | s, a) = P(s' | s, a)$
 if $a' = \operatorname{argmax}_b \phi(\bar{s}, b) \theta_n$
 $\neq 0$ otherwise.

Finally, you can see that this expression again does not involve S prime and hence if you take this and multiply it with this, one can see that we would end up with an expression that is shown over here. And again I would like to highlight that because of the dimensions that are mentioned over here, this will be a D cross 1 vector. Similarly, this will be a D cross 1 vector and this will be a D cross 1 vector. In particular, this matrix that is present, this will be a D cross D matrix and this will be a D cross D matrix. Is this okay?

All right. And given this expression that we saw, it is not difficult to see that the conditional expectation of this expression can be written as B_n minus A_n theta n , where B_n is given by this expression and A_n is given by this matrix. And consequently, we can say that the theta n update rule can be succinctly written as theta n plus 1 is theta n plus some step size times b_n minus a_n theta n plus 1 , where b_n plus 1 is the true update minus the conditional expectation. And in the previous slide, we saw that the conditional expectation is exactly equal to this. So one can see that this exactly equals the true update that we have been studying so far.

<p>Hence,</p> $E[F_n \mathcal{F}_n] = b_n - A_n \theta_n,$ <p>where</p> $b_n = \phi^T D_{\pi_n}^\epsilon \mu \quad f$ $A_n = \Phi^T D_{\pi_n}^\epsilon (I - \gamma P_{\pi_n}^\epsilon) \Phi$	<p>Hence,</p> $\theta_{n+1} = \theta_n + \alpha_n [b_n - A_n \theta_n + M_{n+1}]$ <p>where</p> $M_{n+1} = \delta_n \phi(s_n, a_n) - E[\delta_n \phi(s_n, a_n) \mathcal{F}_n]$
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<p>Hence,</p> $E[F_n \mathcal{F}_n] = b_n - A_n \theta_n,$ <p>where</p> $b_n = \phi^T D_{\pi_n}^\epsilon \mu \quad f$ $A_n = \Phi^T D_{\pi_n}^\epsilon (I - \gamma P_{\pi_n}^\epsilon) \Phi$	<p>Hence,</p> $\theta_{n+1} = \theta_n + \alpha_n [b_n - A_n \theta_n + M_{n+1}]$ <p>where</p> $M_{n+1} = \delta_n \phi(s_n, a_n) - E[\delta_n \phi(s_n, a_n) \mathcal{F}_n]$
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

So, this is what the Q-learning with linear function approximation and epsilon-greedy exploration looks like. In a formal sense, this is what it looks like. In particular, the thing that drives it has the nature that is given over here.

$$E[F_n] = b_n - A_n \theta_n$$

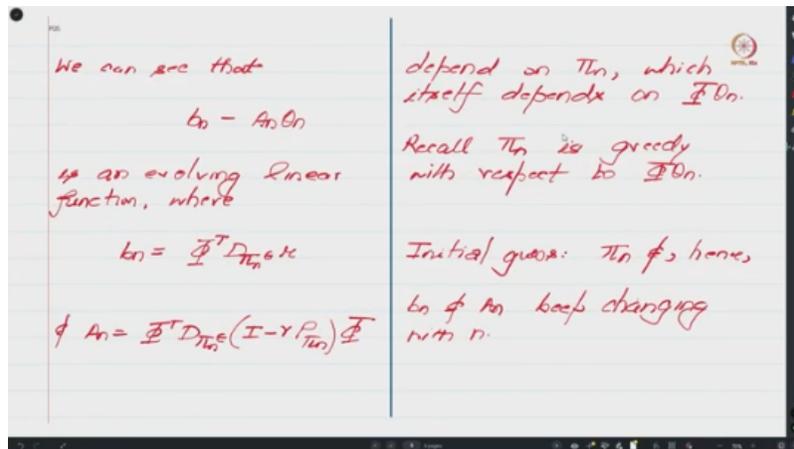
$$b_n = \Phi^T D_{\pi_n^\epsilon} \mu$$

$$A_n = \Phi^T D_{\pi_n^\epsilon} (I - \gamma P_{\pi_n^\epsilon}) \Phi$$

$$\theta_{n+1} = \theta_n + \alpha_n [b_n - A_n \theta_n + M_{n+1}]$$

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

So, at first glance, one can see that this B_n and A_n depend on n . By depend on n , I mean that, as your θ_n changes, your ϕ_n changes.



And because your ϕ_n changes, your π_n also changes. And because your π_n is changing, $\pi_n \epsilon$ will also change. And consequently, your b_n and a_n actually depend on the value of ϕ_n . So, while this may appear to be linear in nature, the challenge in the analysis stems from the fact that this B_n and A_n actually depend on ϕ_n , which means that as your θ_n changes, B_n and A_n also change with them. And this is what sort of complicates the analysis.

And I would like to again highlight that this π_n is greedy with respect to ϕ_n . So, as I said, the summary of this whole discussion is that Because of these reasons, one can guess that because π_n and hence b_n and a_n are changing, the analysis of this algorithm will be very, very difficult.

$$b_n = A_n \theta_n$$

$$b_n = \Phi^T D_{\Pi_n} \epsilon$$

$$A_n = \Phi^T D_{\Pi_n} (I - \gamma P_{\Pi_n}) \Phi$$

And now we will see how to tackle this. So towards that, what we will do is we will introduce a few notations. The first of these notations is this a bar notation.

For $\bar{a}: S \rightarrow \mathcal{A}$, let $R_{\bar{a}} = \{ \Phi \theta : \arg \max_a \Phi^T(R_{\bar{a}})_a \theta = \bar{a}(s) \}$
 Claim: $R_{\bar{a}}$ is a cone
 Proof: $\Phi \theta \in R_{\bar{a}}$ implies $c \Phi \theta \in R_{\bar{a}}$ for any $c > 0$.

Let $b_{\bar{a}} = \Phi^T D_{\bar{a}} \kappa$
 $A_{\bar{a}} = \Phi^T D_{\bar{a}}^2 (I - \gamma P_{\bar{a}}) \Phi$
 Then $b_{\bar{a}} - A_{\bar{a}} \theta \in R_{\bar{a}}$
 $= \sum_s (b_{\bar{a}} - A_{\bar{a}} \theta) \in R_{\bar{a}}$

So, when I write $A_{\bar{a}}$, I am going to imagine this to be a function from the state space to the action space. So, this function takes as input a state and spits out one deterministic action. So, when I say $A_{\bar{a}}$, it is a function which takes as input a state and spits out an action. And what I will do is for any such function \bar{a} , one can also view this \bar{a} as a deterministic policy. Deterministic policy meaning whenever you see s you act according to a \bar{a} of little s . Whenever you see state s you act according to a \bar{a} of little s . So in that sense, this policy whenever it sees a state only picks a particular action and hence I am referring to this \bar{a} as a deterministic policy.

And corresponding to this deterministic policy \bar{a} , one can define a region in the following way. So what is this region? It consists of all those vectors of the form $\Phi \theta$. So θ is d dimensional, ΦA is $S \times A$ cross d , hence $\Phi \theta$ will be a vector that lives in the SA dimensional space where S and A you know are capital S and capital A and hence when I say SA dimensional space it lives in the space that corresponds to the product of the state and action spaces.

So this lives in a very high dimensional space. So it is the collection of all those vectors in this high dimensional space whose greedy policy is actually \bar{a} . So, I should perhaps emphasize this should be true for all S in capital S . So, what I am saying over here is $R_{\bar{a}}$ is the collection of all $\Phi \theta$ s whose greedy policy is \bar{a} .

More formally, if you look at $\Phi^T R_{\bar{a}} \theta$, you look at this inner product and look at the action that maximizes this inner product for a fixed state S . So, you are collecting all those $\Phi \theta$ s where when you try to do this maximization, you end up

with $A \bar{S}$. So and if this happens for all S you put that $\phi \theta$ in this region and $RA \bar{}$ is made up of such $\phi \theta$ s. We will refer to this $RA \bar{}$ because of this nature as the greedy region that is associated with your $A \bar{}$ policy. So given a policy deterministic policy $A \bar{}$ we will refer to $RA \bar{}$ as the greedy region associated with your policy $A \bar{}$. Now our claim is that this region is a cone, right?

That is our claim. So what do we mean by a cone? By a cone we mean that if $RA \bar{}$ contains a vector, then if you take any scalar multiple of that vector, in particular a positive scalar multiple, not arbitrary scalar. So if you take a positive scalar multiple of this vector, then that resultant vector is also lying in that region, okay? That is the interpretation of cone that I am going to use.

So, the claim is that for any deterministic policy $A \bar{}$, your $RA \bar{}$ is actually a cone. So, let us just verify that. So, suppose $\phi \theta$ lies in $RA \bar{}$. Now, $\phi \theta$ lies in $RA \bar{}$ implies that $A \bar{}$ is greedy with respect to $\phi \theta$. In other words, this condition holds true for all little s in capital S .

And now if I multiply this vector $\phi \theta$ by this positive constant c , then one can see that the ordering of the different coordinates of $\phi \theta$, that ordering will not change. Of course the value will change, but the ordering will not change. By ordering I mean that if there was some coordinate of $\phi \theta$, that was larger than some other coordinate, then that same coordinate will continue to be larger than the previous coordinate even after multiplying by c . This is because since c is positive, so I mean the broad idea I am saying is that if you have two real numbers x and y and if x is bigger than y and if I multiply both sides by c , then cx will continue to be bigger than cy . That is the idea that I am using to conclude that

If you have $\phi \theta$ in $RA \bar{}$, then $C \phi \theta$ is not, since you are multiplying each coordinate of $\phi \theta$ by C and C is positive, right, the ordering of the coordinates of $\phi \theta$ is not going to change. Consequently, If $\phi \theta$ is in $RA \bar{}$ which implies $A \bar{}$ is greedy with respect to $\phi \theta$, one can conclude that $A \bar{}$ continues to be greedy even with respect to $C \phi \theta$ which implies in turn that $C \phi \theta$ will belong to $RA \bar{}$ for any C bigger than or equal to 0, strictly bigger than 0. And from this, in fact one

can show it is greater than or equal to 0 as well. And from this one can see that your ϕ_{θ} lying in $R_{\bar{A}}$ implies $C \phi_{\theta}$ lies in $R_{\bar{A}}$ for every C bigger than or equal to 0 and hence one can conclude that $R_{\bar{A}}$ is a cone.

So the summary is that if you define a greedy policy in this way, And if you look at the, sorry, I should say it again. If you define a deterministic policy in this way and you define the greedy region associated with this deterministic policy in this way, then this greedy region one can immediately show is actually a cone. So now define $B_{\bar{A}}$ to be this vector over here. So this vector is very similar to your B_N but the difference is that

Here instead of having $\pi_N \epsilon$ we have a $\bar{\epsilon}$. So what is a $\bar{\epsilon}$? It is the epsilon randomization of a \bar{A} which means \bar{A} is a deterministic policy with $1 - \epsilon$ probability you act according to \bar{A} and with ϵ probability you pick a random action. That is what $\bar{A} \epsilon$ is.

So this is again a stochastic policy. You look at the stationary distribution that is associated with this stochastic policy. And then you know define this matrix over here. Again this is a $S \times S$ matrix whose S_{th} entry is basically the product of the stationary distribution at state S associated with this policy and the probability of the action under this stochastic policy. That is what this matrix would be and accordingly one can see that $B_{\bar{A}}$ is defined in this way.

And similarly, capital \bar{A} is defined in this way. Again, this has a very close structure or similarity to the A_N matrix that we saw before, except that there, you know, this quantity and this quantity were $\pi_N \epsilon$ and π_N respectively. Here instead, we work with, I should perhaps clarify that this is $\bar{A} \epsilon$. So, here instead of having $\pi_N \epsilon$ and π_N , we have $\bar{A} \epsilon$ and \bar{A} . So, in this way you define this $B_{\bar{A}}$ and \bar{A} .

Then one can see that $B_N - A_N \theta_N$ can actually be written in the form that is given over here. In particular, one can see that $B_N - A_N \theta_N$ can be expressed as a sum of $B_{\bar{A}} - \bar{A} \theta_N$. So, $B_N - A_N \theta_N$ can be expressed as sum over \bar{A} . $B_{\bar{A}} - \bar{A} \theta_N$ along with an indicator, so what is this indicator? The indicator over here is ϕ_{θ_N} belongs to $R_{\bar{A}}$, in other words.

Whenever your theta N is such that it belongs to R A bar, right? Then one can see that the B N minus A N theta N will actually be B A bar minus A A bar theta N, okay?

For $\bar{a}: S \rightarrow \mathcal{A}$, let

$$R_{\bar{a}} = \{ \bar{\theta} : \forall s \in S, \arg \max_a \phi^T(R_{\bar{a}}) \theta = \bar{a}(s) \}$$

Claim: $R_{\bar{a}}$ is a cone

Proof: $\bar{\theta} \in R_{\bar{a}}$
 implies $c \bar{\theta} \in R_{\bar{a}}$ for any $c > 0$.

let $b_{\bar{a}} = \bar{\Phi}^T D_{\bar{a}}^{-1} \kappa$

$$A_{\bar{a}} = \bar{\Phi}^T D_{\bar{a}}^{-1} (I - \gamma P_{\bar{a}}) \bar{\Phi}$$

then $\sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{ \theta_n \in R_{\bar{a}} \}}$

$$= \sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{ \theta_n \in R_{\bar{a}} \}}$$

So, I would like to first emphasize a bit why we can write it in this following way, right? So, for every deterministic policy A bar, we can come up with a greedy region associated with A bar. And one can show that if you have two different greedy regions except for the origin. If you sort of decide some rule to, you know, break ties and so on, one can see that except for the origin, okay, these two regions will actually be disjoint, right? So, except for the origin, these two regions will actually be disjoint, right?

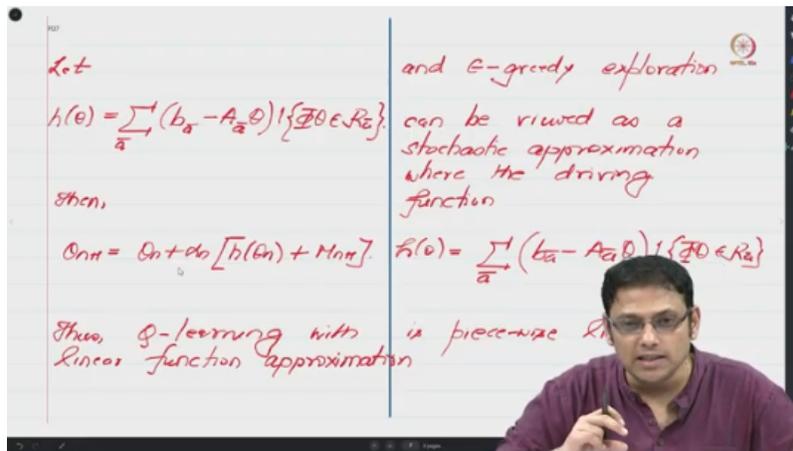
So, in other words, the greedy regions will be pairwise disjoint and one can show that the whole RSA space, okay, the whole RSA space can be written as a union of these greedy regions, right? So, whenever you have a phi theta n which is not equal to the origin, one can show that phi theta n must lie in R A bar for some A bar, right? So, accordingly, and in fact, one can show that when phi theta n is not the origin, right? one can show that there is exactly one A bar, right? So there is one A bar and that A bar is unique such that phi theta n belongs to R A bar. Again, I am presuming there is some way to break ties and, you know, if you sort of think a bit closely, you will be able to see why I require this breaking of ties and so on, right? But once you have a rule to break ties, one can conclude that phi theta n

<p>For $\bar{a}: S \rightarrow \mathcal{A}$, let</p> $R_{\bar{a}} = \left\{ \Phi^T \theta : \forall s \in S \right. \\ \left. \arg \max_a \phi^T(R_a a) \theta = \bar{a}(s) \right\}$ <p>Claim: $R_{\bar{a}}$ is a cone</p> <p>Proof: $\Phi^T \theta \in R_{\bar{a}}$ implies $c \Phi^T \theta \in R_{\bar{a}}$ for any $c > 0$.</p>	<p>let $b_{\bar{a}} = \Phi^T D_{\bar{a}} \kappa$</p> $A_{\bar{a}} = \Phi^T D_{\bar{a}}^2 (I - \gamma P_{\bar{a}}) \Phi$ <p>Then $\sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{\Phi^T \theta_n \in R_{\bar{a}}\}}$</p> $b_{\bar{a}} - A_{\bar{a}} \theta_n$ $= \sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{\Phi^T \theta_n \in R_{\bar{a}}\}}$ <p>For $\bar{a}, \bar{a}' \quad R_{\bar{a}} \cap R_{\bar{a}'} = \{0\}$</p>
-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<p>For $\bar{a}: S \rightarrow \mathcal{A}$, let</p> $R_{\bar{a}} = \left\{ \Phi^T \theta : \forall s \in S \right. \\ \left. \arg \max_a \phi^T(R_a a) \theta = \bar{a}(s) \right\}$ <p>Claim: $R_{\bar{a}}$ is a cone</p> <p>Proof: $\Phi^T \theta \in R_{\bar{a}}$ implies $c \Phi^T \theta \in R_{\bar{a}}$ for any $c > 0$.</p>	<p>let $b_{\bar{a}} = \Phi^T D_{\bar{a}} \kappa$</p> $A_{\bar{a}} = \Phi^T D_{\bar{a}}^2 (I - \gamma P_{\bar{a}}) \Phi$ <p>Then $\sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{\Phi^T \theta_n \in R_{\bar{a}}\}}$</p> $b_{\bar{a}} - A_{\bar{a}} \theta_n$ $= \sum_{\bar{a}} (b_{\bar{a}} - A_{\bar{a}} \theta_n) \mathbb{1}_{\{\Phi^T \theta_n \in R_{\bar{a}}\}}$ <p>For $\bar{a}, \bar{a}' \quad R_{\bar{a}} \cap R_{\bar{a}'} = \{0\}$</p> $\mathbb{1}_{\{ \cdot \}} = \bigvee_{\bar{a}} R_{\bar{a}}$
-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

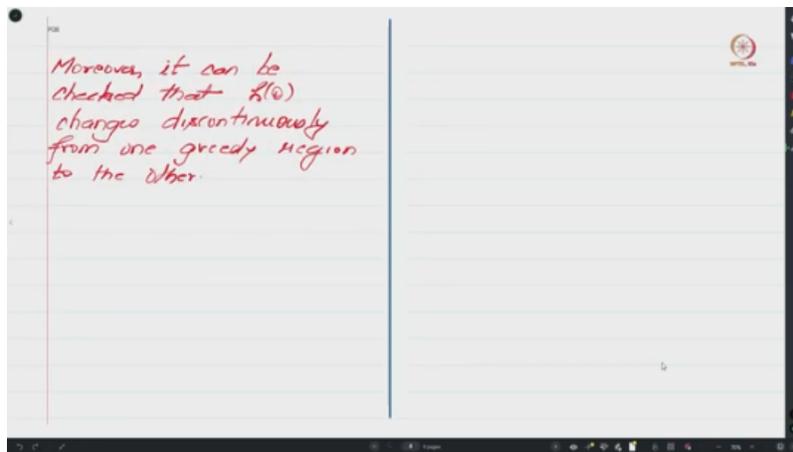
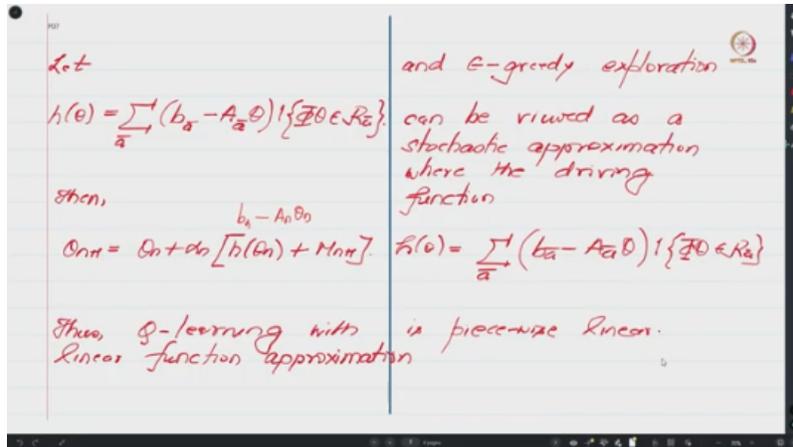
So for every theta, or every theta n, there should be a unique A bar such that phi theta n lies in that R A bar. In other words, this indicator, whenever phi theta n is not 0, will be 1 for some A bar, and you know there will be a unique A bar for which this happens. In other words, we can write bn minus An theta n in the following way. Now, this is very interesting because when you have finite state and action spaces, one can show that the number of deterministic policies is finite, which means the number of greedy regions is also finite, and one can think of this as a partition of this space into the different greedy regions. And what this says is that when you are, you know, phi theta n or theta n lies in a particular greedy region, right, or phi theta n lies in some particular greedy region, right, then the dynamics Bn minus An theta n can be expressed in the way that is shown over here.

So accordingly, what we are going to do is we are going to define h of θ in the following way. This is sum over A bar of b A bar minus A A bar θ indicator ϕ θ belongs to R A bar. Then one can see that the update rule for Q-learning with linear function approximation can be written as θ_{n+1} equals θ_n plus α_n times h of θ_n plus m_{n+1} . So before we had this expression b_n minus A θ_n , right?



So, when we looked at that expression, it appeared that your b_n and A_n keep changing with respect to every n , right? However, what this description of s and this update rule that we have written over here suggests that you know, when your ϕ θ_n lies in a particular greedy region, then you follow this linear dynamics. And when your ϕ θ_n moves from one greedy region to the other, then the linear dynamics also changes. Is this okay?

So this is the key observation that we need to make. In other words, one can note that or formally one can see that this h of θ that I have defined over here is actually piecewise linear. That is, in one greedy region it has some linear nature, and when you move from one greedy region to the other, then the dynamics also changes accordingly. And then, the surprising or the challenging aspect of this definition of H is that this H actually changes discontinuously when we move from one greedy region to the other. So, this brings us to the end of this class.



Let me quickly summarize what we have done over here. So, in this class, we tried reading Q-learning with linear function approximation and epsilon-greedy exploration. And we wanted to analyze the behavior of this algorithm. In particular, in the simulations, we saw some strange behavior, right? And we wanted to understand: can we say something about, you know, why this asymptotic behavior is happening? Can we explain something towards that?

So, in that direction, what we did was we tried, you know, taking the conditional expectation, which is the standard thing we do, of the update and, you know, adding and subtracting them, and we saw that the update rule is of the form b_n minus $a_n \theta_n$, right? And from the definition of b_n and a_n , we saw that To begin with, it appeared that b_n and a_n both change with your iteration n . In particular, they depend on $\phi \theta_n$. So, as your $\phi \theta_n$ changes, your b_n and a_n also change. But then we

observe that if we define these greedy regions associated with these different deterministic policies, then we observe that this $b_n a_n \theta_n$, this dynamics that governs this Q-learning with linear function approximation, does not change arbitrarily. The behavior is that whenever you are in a particular greedy region, it is the linear dynamics that is B_A bar minus A_A bar theta associated with that greedy region that governs your dynamics. And when you move from one greedy region to the other, there is a different dynamics that governs your behavior.

And in the next class, we will see how this changing dynamics from one greedy region to the other affects the behavior of the algorithm, and we will see how this discontinuously changing, piecewise linear function can be used to describe the different behaviors that we saw in those empirical simulations. So, until then, goodbye and namaste, and thank you.