**Advanced Numerical Analysis**
**Prof. Sachin Patwardhan**
**Department of Chemical Engineering**
**Indian Institute of Technology – Bombay**

**Lecture – 38**
**Solving Nonlinear Algebraic Equations: Introduction to Convergence Analysis of**
**Iterative Solution Techniques**

We have been looking at nonlinear algebraic equations and we looked at three different classes of methods. One was derivative free method, the other was sloper derivative based methods and the third was optimization, so which was numerical optimization and we are looked at the algorithmic aspect of nonlinear algebraic equations. Now, today I am going to touch upon the convergence aspect.

So very, very important aspect of equations, but I am just going to give a very, very brief introduction. I am not going to go deep into this. I just want to sensitize you that there exists lot of work, lot of literature on convergence of nonlinear iterative schemes. For convergence of linear iterative schemes like Gauss-Seidel method, Jacobi method, we could actually derive in the class necessary and sufficient conditions.

Whereas for nonlinear cases, much more difficult and the machinery that you require it is fairly more advance than what we are doing covering in this course and also many times, you only get sufficient conditions, you do not necessary conditions. So, nevertheless these tools or the theorems that actually give sufficient conditions give lot of insight into how solutions of nonlinear algebraic equations behave.
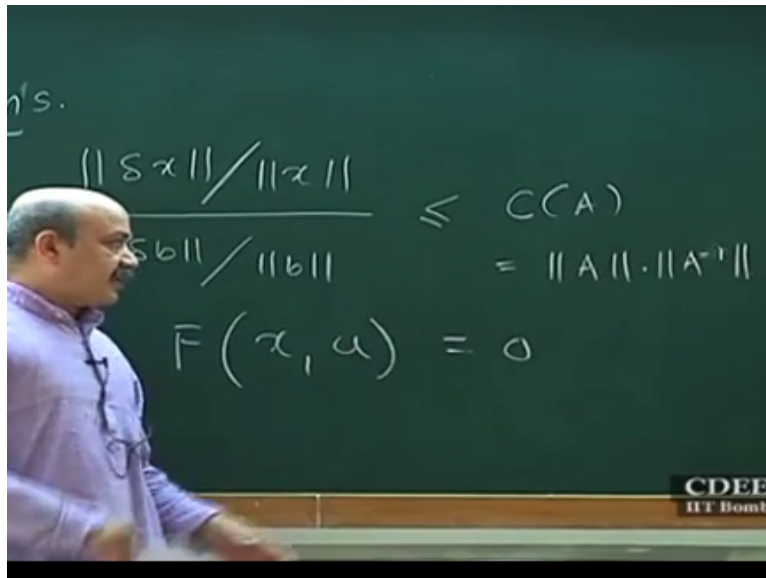
So, I am just going to touch up on it today, not really go into deep of this subject. So, one thing that we need to talk about, see if you look at the development that we did for the linear algebraic equations. We had there is some sense parallel between what we have done there and what we done here. There too we talked about noniterative schemes, iterative schemes and then we talked about optimization based schemes.

There we talked about a very important issue called condition number. So, we said that a set of linear algebraic equations is well conditioned or ill conditioned depending upon some properties of matrix A, right and it was possible to do analytical treatment quite easily with

whatever we have learnt till now. Can we extend this to nonlinear algebraic equations, I am just going to briefly touch up on this idea.

That what is the condition number of a nonlinear algebraic system and then move onto the convergence properties of, or how do we analyze the convergence of nonlinear algebraic equations? So, what was in the case of, so first thing I just want to touch up on this condition number. So in linear algebraic equations, we had defined condition number as when you have $Ax = b$, okay.
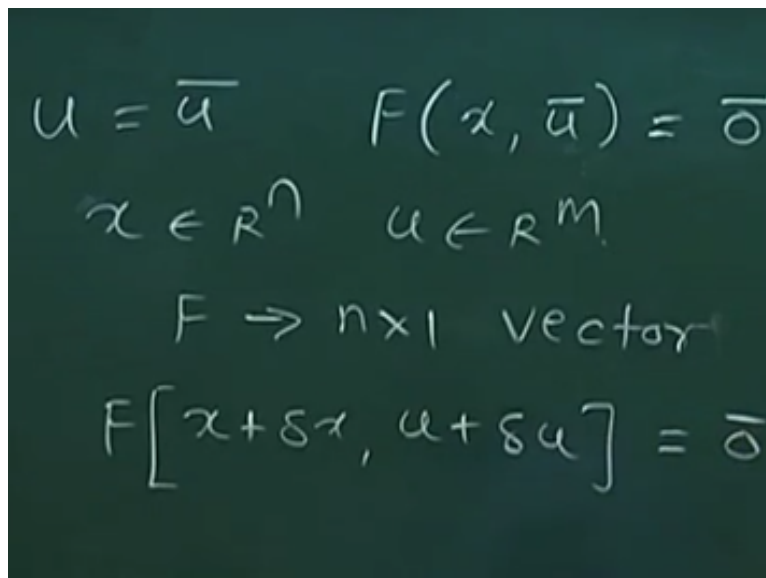
One may be defined this condition number was sensitivity of the solution x to a small change in b, right. So, if we look at this as an input and x as an output, if you look at this first time as b as a input and x as output where A is the operator, one maybe defined the condition number was norm delta x/norm x. So, we showed that this ratio that is fractional change in the solution to fractional change in the input, okay is bounded by this condition number, which is multiplication of norm of A * norm of A inverse.

Now to draw a parallel, I am going to consider nonlinear algebraic equations of the form f of x u = 0, okay. Well, this kind of equation very routinely arise in chemical engineering when you are solving steady state behavior of say CSTR, x are no states concentration temperature inside the reactor, u are inputs, as an input flow rate, inlet concentration, inlet temperature, all these free parameters, okay, input parameters.

So if you fix yourself to one input condition, you will get one steady state of the reactor or let us say you have distillation column, you have this kind of equation, u there is nothing but feed composition, feed flow rate, feed temperature, reflux rate, heat input, okay all these free inputs which in balance of control you call as disturbances or manipulate inputs, all these inputs are u.

X are all the dependent variables like tray temperature, tray concentrations, vapor concentration, liquid concentration everything. So moment to fix one u vector, okay, for a particular u vector, let us say u = u bar, you get f of x u bar = 0, this is what you have to solve. Once you fix u bar, okay, say typically x belongs Rn and u belongs to Rm and for every u, you fix, okay and f is a n cross 1 vector, okay.

**(Refer Slide Time: 06:46)**



This particular vector, n cross 1 vector is a state of nonlinear algebraic equations. You have to solve them simultaneously for a given u, for every u, okay. If I change the feed condition, if I change the feed composition, the concentration or temperature profile in the distillation column is going to be different, okay. For every value of this input conditions, you get one set of steady state solution, okay.

In some sense, this u is parallel with b on the right hand side. If you change the right hand side b, you get different x, okay. So now, we define condition number as with respect to solution of f of x + delta x and u + delta u = 0. So, when I change u for u bar 2, say u bar + delta u bar, okay. When I introduce a perturbation in u, what is the corresponding perturbation in the solution x, okay.

So, I am going to define sensitivity of this equation or sensitivity of the solution with reference to perturbation in u as my condition number, same idea fractional sensitivity of the solution to fractional sensitivity or fractional change in the input, okay that is going to be my condition number. So for a nonlinear system, we define cx as supremum over delta u, okay. We define this a supremum over all perturbations delta u, okay.

**(Refer Slide Time: 09:04)**



This of course nonzero perturbations delta u. So in other words this delta x, not c of x, c of f, this should be, well this in general will not be a constant number like matrix, you get you know matrix is a operator which only consist of, we have considered a matrices which are of real numbers, so you will get you know matrix, norm of matrix * norm of matrix inverse as your condition number.

Here that is not going to happen, your nonlinear algebraic equations. So, the conditioning of a nonlinear algebraic system could be different in different regions of the state space. Suppose, you are solving, this is a abstract way of putting it, I will put it in the simple words, let us say distillation column, you are trying to solve set of algebraic equation for a binary distillation column in a low purity region as against in the high purity region, okay.

Conditioning of these nonlinear algebraic equations in low purity region, okay will be different from conditioning of these nonlinear algebraic equations in the high purity region. It might be more difficult to solve, for example high purity region. I am not saying it is always

difficult, but little bit it might be more ill condition let us say and it is well condition when you are away from the high purity region.

So if you are trying to solve equations when the purity is you know 0.99 as against purity to top purity as against top purity is 0.9, you will have different behavior of the nonlinear algebraic equations, okay. So, the sensitivity of the solution to a small change, okay on the right hand side might be different in different regions. It depends upon where in the state space, you are solving this set of equations that is critical, okay.
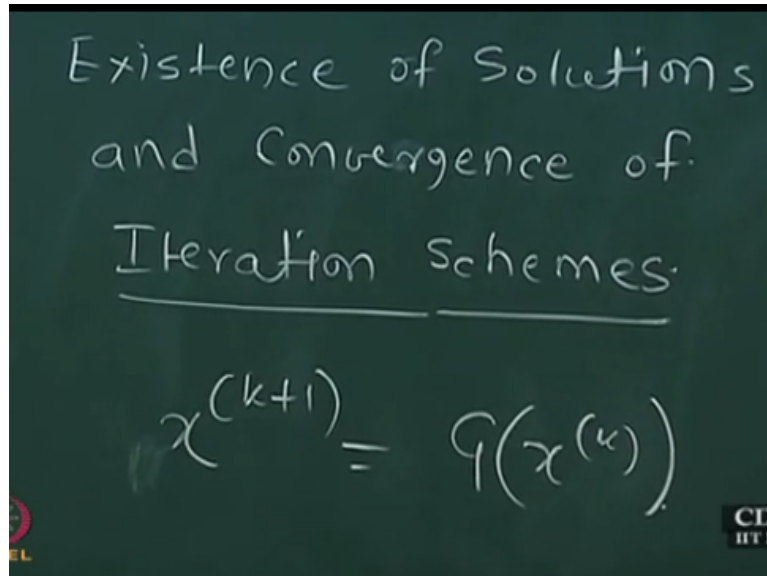
So, this actually gives you an upper bound on this ratio change in the or a perturbation fractional change in the solution to fractional change in the input condition, okay. So, analogous to the linear case, one can define something called a condition number here, you can talk about well conditioned nonlinear systems, ill conditioned nonlinear systems, you can talk about local, you have to understand this is local, okay.

For a same nonlinear system, it could be well condition in some region, it could be ill condition in some region, okay. So, nonlinear algebraic equations are much more difficult to handle in terms of conditioning than linear algebraic systems. So sensitivity, what did actually condition number tell you? Sensitivity of the solution to errors for example, okay.

So, if nonlinear algebraic equations are you know in some cases, if the condition number is high, which means the small change in u will cause a large change in the solution x, okay. A small error in representation of u will cause a large change in the solution x and just imagine when we are solving many of these nonlinear algebraic equations arise because we are doing discretization of some nonlinear boundary value problem or some partial difference equations.

When you doing that you are approximating, okay, so in some regions, a small perturbation in the input condition can lead to a large change in the solution because of sensitivity of the equations in that region. But this is again, as I said it is much more difficult to analyze this than the linear case. The next concept is we just touch up on this existence of solution and convergence of iteration schemes.
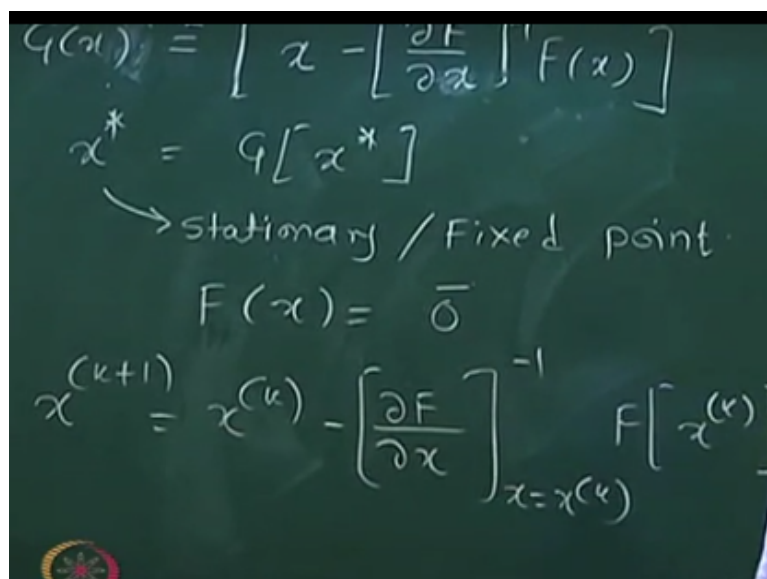
**(Refer Slide Time: 14:34)**

Existence of Solutions and Convergence of Iteration Schemes.

$$x^{(k+1)} = G\left(x^{(k)}\right)$$

Now, you have seen that all the methods that we have for solving nonlinear algebraic equations or iterative. Quadratic multidimensional equations can be solved analytically, but I am not aware of solution for the cubic case, so majority of elements in the set of nonlinear algebraic equations cannot be solved analytically, you have to solve them using some numerical procedure, okay.

Invariably any numerical scheme that you come up with can be written in this form, any numerical scheme that you come up with, okay. You start with a guess, generate a new guess, x =, where you want to reach finally, I want to reach finally to what is called the stationary point, I want to reach to a stationary point x star = G of x star. This x star is called as stationary point, it is called as fixed point, okay.

**(Refer Slide Time: 16:17)**



$$G(x) = \left[ x - \left[ \frac{\partial F}{\partial x} \right]^{-1} F(x) \right]$$

$$x^* = G\left[ x^* \right]$$

$$\rightarrow \text{stationary} / \text{Fixed point}.$$

$$F(x) = \bar{0}$$

$$x^{(k+1)} = x^{(k)} - \left[ \frac{\partial F}{\partial x} \right]^{-1}_{x = x^{(k)}} F\left[ x^{(k)} \right]$$

So, we want to actually reach here. See for example, when you are solving f of x = 0, if you are solving using Newton-Raphson method, okay or Newton's method, Newton's method was x k+1 = xk - doh f/doh x at x = xk inverse f of, right this was my Newton's method. I wanted to solve for f of x = 0. Now, G would be here, in this case, G is equivalent to x - doh f/doh x inverse fx, this is my G of x, okay.

And ultimately you are solving for x = G of x, right. You are solving for x k+1 = G of xk. From the previous guess, you construct a new guess, okay. So any method that we are looked at till now for solving nonlinear algebraic equations, iterative method can be put into this generic form and you are looking for x star, x star is the fixed point, okay. I think the word stationary is not really used here mostly.

Stationary point is used in the case of optimization; it is the fixed point. So, literature on function analysis will be full of fix point theorems, so doing analysis of iterative equations, okay, so how do the, okay. So now, I am going to just revisit some other terms that we looked at right in the beginning Banach space and operator mapping, Banach space to Banach space and so on, okay.

Why I am worried about Banach space? What is the Banach space? Banach space is one in which every sequence has a limit within the space is convergence. Why I am worried about every sequence, looked here. What is this? If I am start from some x not, okay, I will get a sequence of vectors x1, x2, x3, x4, x5 and so on. This iterative process, we generate a sequence of vectors, right.
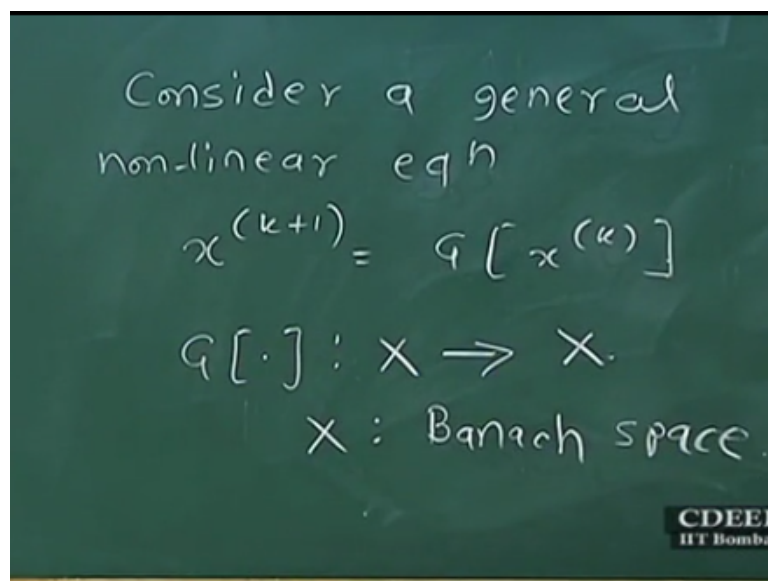
Now if I give one particular problem, okay and if I ask him to solve the problem, he will start with one x not, she will start with another x not, she will start with another x not, okay. What is important is that if they are starting from different initial guesses, okay, will those sequences converge to the same solution under what condition. First of all, one condition or one primary condition is that the sequence should not go to a limit which is outside the space, right.

The sequence should remain within the space that is the first condition. Second condition that is important is that we want to know is that whether the sequence will converge to a solution,

is the solution unique? So, does the solution exist and is the solution that you get is that unique, all these questions are very, very important, okay.
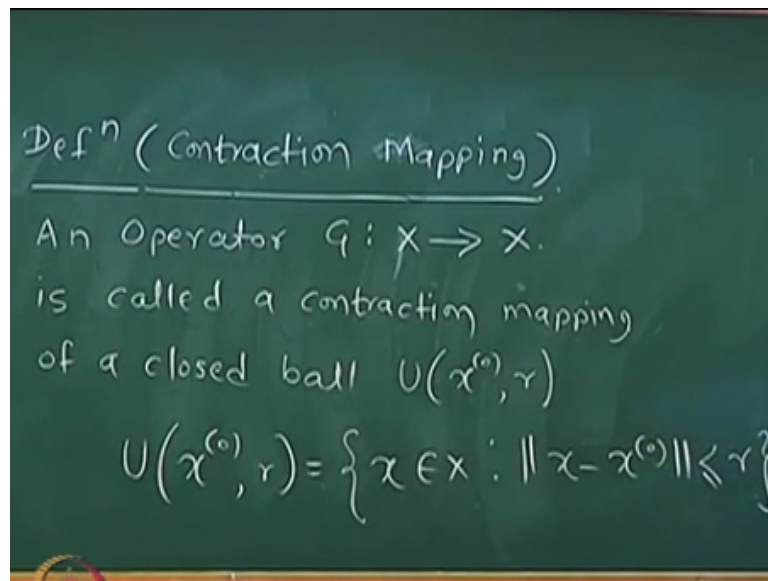
So, I am just going to give some hint about how these are handled in the, so in some sense this would connect to the cherry that we had done in the beginning, you know abstractions of Banach space and Hilbert space and so on. Now, G is the mapping from x to x, where x is the Banach space or a complete normed linear space which means moment I say this I am ensuring that the sequence generated from any initial guess x not will never leave the space, will always be within the space that is what I mean here, okay.

**(Refer Slide Time: 21:16)**



The sequence will never leave the space. An important concept here is contraction mapping, okay. A very important concept here is contraction mapping. Now, when I am writing here an operator G implicitly one which was define, we have just define it is also x is Banach space to Banach space, all these things are implicit, I am not writing them on the board, okay. I will just, I have to complete this definition, but before that let us look at what I have written here.

**(Refer Slide Time: 22:25)**

$$Def^n \ (Contraction \ Mapping)$$

An Operator $G: X \Rightarrow X$.
is called a contraction mapping
of a closed ball $U(x^{(o)}, r)$

$$U(x^{(o)}, r) = \left\{ x \in X : \| x - x^{(o)} \| \leqslant r \right\}$$

An operator G is called as a contraction mapping of the closed ball, okay. A closed ball is set of all x belonging to the vector space x, such that x - x not is < r, r is some radius, okay. How do you, what is the relation of this radius and convergence all that will come to soon, but right now, I am defining idea of contraction, okay on a small vision in the neighborhood of x not.

This is the way of defining the neighborhood of x not, some region around x not, okay. So, which norm you use? Depends upon you, 1 norm, 2 norm, infinite norm, it does not matter, any norm that of your choice, but I am defining a region in the neighborhood of a initial guess, okay. What is x not here because we are solving nonlinear algebraic equations we can look at x not as my initial guess, okay.

It is not a fixed point, as I said x not can vary from person to person, everyone can take a different guess, okay. Just pay attention to these concepts because these are little difficult and then you are not, the other things which have been teaching at least you know something about it, okay. Whereas these are little advance concept, so you have to understand them carefully.

**(Refer Slide Time: 26:08)**

if there exists a real number
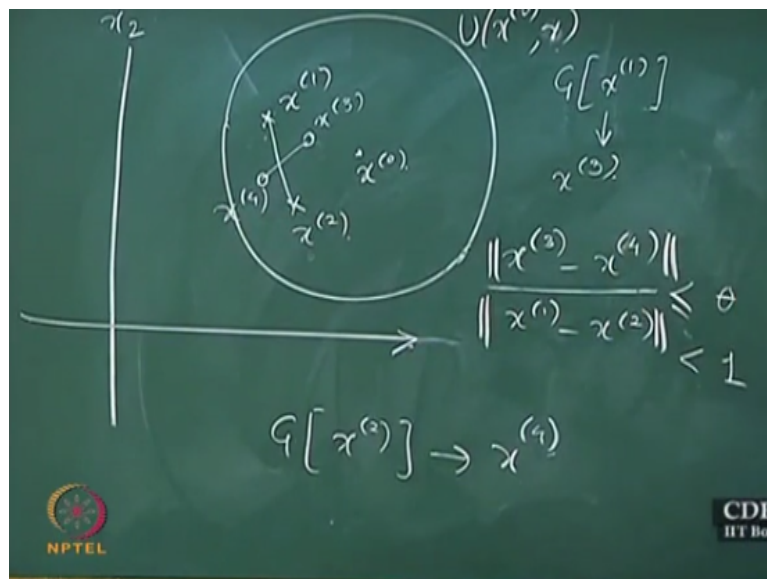$0 \leq \theta < 1$ such that
$$\| g[x^{(1)}] - g[x^{(2)}] \|$$
$$\leq \theta \| x^{(1)} - x^{(2)} \|$$
for all $x^{(1)}, x^{(2)} \in U(x^{(0)}, r)$

Now, I want to call this mapping into a contraction mapping, if there exists a real number theta which is strictly < 1, which is a positive number strictly < 1 such that, okay. So, this completes my definition. So when do I call mapping G to be a contraction mapping, okay. If I pick up any two points x1 and x2 in this region, okay and take difference between G of x1 and G of x2 that is always smaller than x1 - x2, which means if I draw it pictorially.

**(Refer Slide Time: 28:09)**



Let us say this is my x2 and x1 and this is my x not, initial guess and let us say this is the region, this is the ball in which I am defining the contraction mapping, okay. What I am going to do is I am going to randomly pick any two points, say here and here, okay. Now, what is G? G is an operator, which gives you element in the same set, right. G is the mapping from x to x.

So, if I apply G on one element, it will give you another element, okay. So, let us say this is my x1 and this is x2, okay. So, when I apply Gx1, see what is G? x = G of x, right. This is the kind of equation we are solving, so when we apply G on x, you get another x, okay. So, let us say this gives me some element x3, okay. I pick up x2 and apply G of x2, this gives me say x4, okay.

Now, we are concerned about this ratio that is x3 - x4 upon x1 - x2. We are saying if this is < theta which is < 1, okay. See I get two points, let us say when I apply this, I get x3 and when I apply G on this, I get x4. What we are saying is that this distance between x1 and x2 is larger than x3 and x4. Sorry, I should put norm here. It is not, we are working in multiple dimensions, I should put norm, okay.

What I am saying here is that the distance between any two points, x1, x2, okay, let us say x1, x2, this distance is always larger than this distance. This is x3 which was obtained by applying G on x1. This is x4 which was obtained by applying G on x2, okay. So, this is my x3, this is my x4, okay. So, if this condition holds for any two x1, x2 inside this region, okay which means when you apply G on x, on any two separate points, okay, then the relative distance contracts.

It comes close and it is called as a contraction map, is this clear, okay. Yeah, **"Professor - student conversation starts"** (()) (32:28) that is the good question, will come to that, okay. "**Professor - student conversation ends".** So that will depend upon how you are chose this radius and it is a very good question, leading question, I will answer this question soon, okay that actually forms the crucial, it is very crucial to the solution procedure, the convergence of solution method.
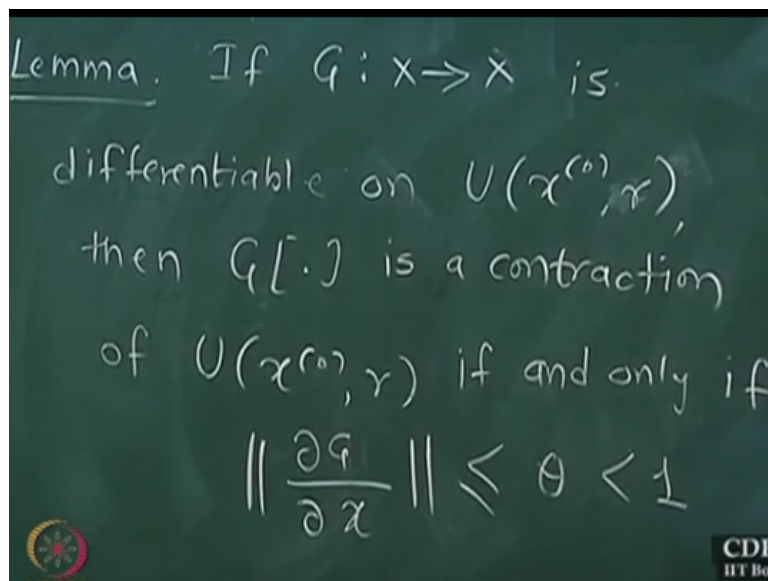
So, let us for the timing assumes that it lies within the same ball, let us assume for the time being. Then, every time you apply on any two points, the new two points that you generate x3 and x4 are closer than the initial two points, you take any two points, apply G on the first point, apply G on the second point, you get two new points, okay that should be closer than x1, x2.

It should happen for any x1, x2, in this region, then G is called as contraction mapping on this ball u, so this is my u x not, r and she has rightly guessed this critical point is, what is this r,

will come to that. Now, in general we are solving for x k+1 = G of xk, it is quite likely that G is not a continuous operator, not a differentiable operator, it could be continuous operator, but not a differentiable operator.
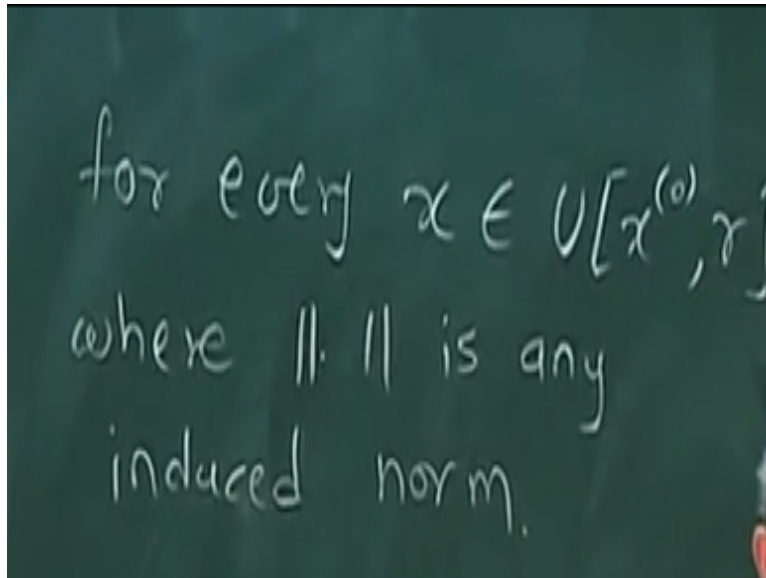
So, actually this theory that has been derived is not necessarily for all differentiable operator, but if G is differentiable which is the case in most of the chemical engineering situations, then we can derive some nice conditions. So this is the result, well all the other things hold that is G is the operator from Banach space to Banach space and it is differentiable on this ball, okay.

**(Refer Slide Time: 34:44)**



Well, this makes it easier for you to understand because derivative is something which you are more comfortable with. So, if the derivative of G is norm of derivative of G is strictly < 1 for every x belonging to, this is the very nice result. It says that if the operator is differentiable, okay, then it is a contraction mapping, if and only if necessary and sufficient condition, if and only if the norm of the derivative is strictly < 1, okay.

**(Refer Slide Time: 36:33)**

for every $x \in U[x^{(0)}, r]$
where $\| \cdot \|$ is any
induced norm.

So, if the norm of the derivative is strictly < 1 in some region, then it is, well, I have to check whether it is necessary and sufficient, I will confirm this. If it is strictly < 1, okay, it is definitely a contraction, but if it is a contraction, does not mean that norm has to be strictly < 1 that we have to check. I am definitely sure that if part of it, I will confirm this result.
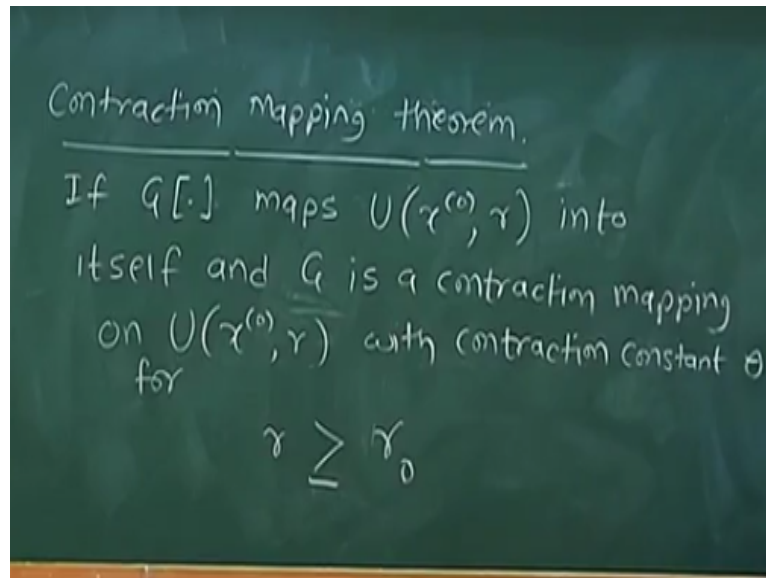
So only if part is in doubt, so if this is strictly < 1, okay. Then, it is surely a contraction. So, if the derivative has norm strictly < 1, we are guaranteed that. So, this part I am not too sure right now, I have to confirm, okay. How are you going to use this contraction mapping business? The literature on theoretical numerical analysis is full of what are called as fixed point theorems.

They are worried about under what condition the solutions to x = G of x exist, under what condition iteration sequences will converge to the solution. The solutions are local, first of all you understand that unlike linear algebraic equations when A is nonsingular, you have a unique solution, right that is not a case in nonlinear algebraic equations. You can have multiple solutions to same set of nonlinear algebraic equations.

Simplest example I have given you is, you know from the abstract this thing is eigenvalues. When we looked at eigenvalue problem, it was a set of nonlinear algebraic equations and any how multiple solutions to that problem. Other example of course is CSTR, a CSTR can have multiple steady states. Under the same input conditions, it can have the steady state operating point and unsteady operating point depending upon how the heat removal and heat generation terms are.

So, same set of nonlinear algebraic equations under identical input conditions can have multiple solutions, okay. So, we are talking about local convergence to local solutions. We are talking about convergence inside a ball, okay, this ball which is in the neighborhood of the initial guess, okay. Now, let us try to under this theorem, this is the contraction mapping principle, one of the fundamental results in.

**(Refer Slide Time: 40:33)**



Now, probably you can already guess norm of an operator strictly < 1, okay. Then, you get convergence. We have seen something similar to this, what was that linear algebraic equations, we were analyzing convergence of iterative schemes and we said that induced norm is a upper bound on the set of you know, the lower bound of that is the spectral radius and so, if norm is < 1, norm of the operator. So, the norm of the operator there was A, okay.

Not A, s inverse t, now the operator there was s inverse t and if now the operator s inverse t was < 1, we were ensured convergence. So, this is something like generalization, so try to compare draw parallels, then you will understand these things better, okay. Now, I am going to assume something which she was suspecting, okay. The theorem assumes that G is the map, which maps you into itself.

So which means you take any point inside this u, okay and apply G on it, the resultant will also be inside u that is the first assumption. So, actually choose r becomes very, very critical because you know G has to map into itself, okay. Now, here I am coming up with the

condition how I choose r, okay. Now, see carefully you have this map G, which is the contraction map, first of all G maps u into itself.

$$r_0 = \frac{1}{1-\theta} \| G[x^{(0)}] - x^{(0)} \|$$

$$= \frac{1}{1-\theta} \| x^{(1)} - x^{(0)} \|$$

Then

(1) G has a unique fixed point $x^*$ in $U(x^{(0)}, r)$

$$x^* = G(x^*)$$

If I take any element in the set u, G will map it into itself. So, you will find a new element also inside u, it is not going to be different. Second thing it is a contraction map, okay which means you take any two points in u and apply G to it, okay. The new point generator are going to be closer than the two initial points, any two points, okay, this is the second thing. What should be the minimum size of this ball? Okay.

Looked carefully it is related to the first x that you produce, okay. **"Professor - student conversation starts"** Why this is related to first x that you produce? (()) (45:52). "**Professor - student conversation ends".** See what should happen is that if you take x1, x2 and x2, x3, okay, x2, x3 will be shorter than, sorry should take x0, x1 and x1, x2 because it is a contraction.

X1, x2 will be shorter than x0 and x1. The very first x1 that you produce by applying, so this is the okay should be > this in some way it is related to distance x1 - x0. How this factor comes, you will have to read the proof. Why just 1 - theta comes, okay, but you can appreciate that the radius is related to the first, if you start with x not, the first x1 that you generate, okay that should be within the ball.
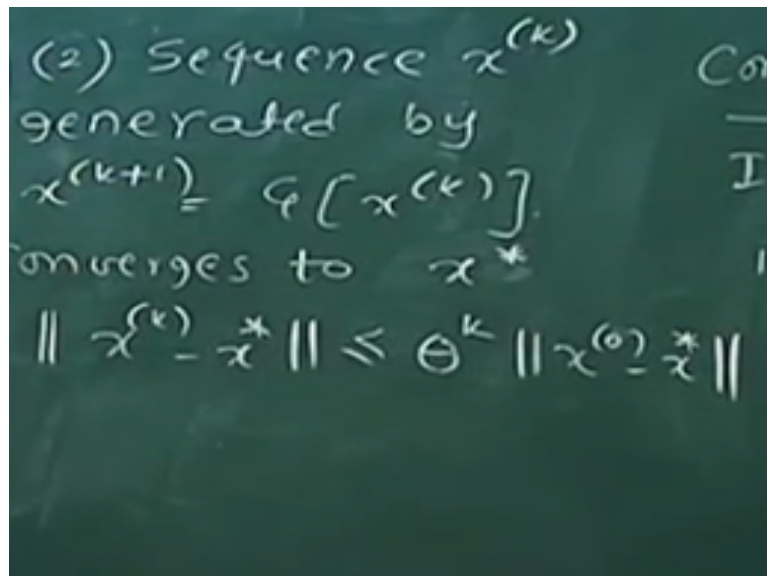
After that whatever you do will be within the ball because it is a contraction, okay. It will stay within the ball, okay. What next? Then, now if these conditions are satisfy, then okay. First

thing that this theorem guarantee is that G has a unique fix point inside the ball, okay. There exists a unique solution inside the ball. What is the solution of the problem? The fix point. You want to reach x star = G of x star, okay.

So, this is the unique fixed point inside this ball, okay. When the radius of the ball is chosen according to this condition, okay, this minimum radius and when G is a contraction on this particular ball, okay, then we are guaranteed that the solution exists inside the ball that exists one point, okay, where this condition is satisfied, okay. Moreover, with this ball is only one such point, there are no two points, okay.

There is only one such point in which a unique solution that is also. Now, the second part is very, very important, I will just continue here. Second and third part, there are three parts for this result. Second part says that if it is a contraction and if these conditions are met, then applying G repeatedly on the sequence will take you to the solution, okay that is guaranteed and at what rate you will go to the solution, okay.
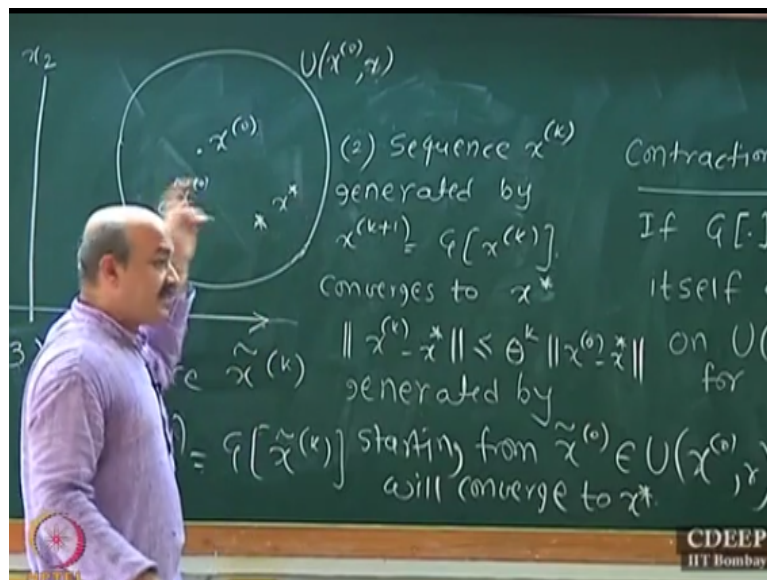
**(Refer Slide Time: 49:09)**



The distance between xk - x star, this will reduce with theta to power k. Again, look at this result, it says that the distance between xk and x star will be shorter than distance between x not and x star. This is the initial distance you started with. X star is let us say this is my x star. This is the solution. I am starting with some x not here, I want to reach here, okay. In doing so, I might move around, you do not around how it will happen.

It is a nonlinear map, you might move around all over the set and then come back to the solution, okay. How the path is going to be, you do not know, but what you know is that the initial distance, okay. Now, how is this result going to shrink rest to theta to power k. Theta is the fraction. So, theta to power k, as k increases, this distance will reduce, okay. If theta as you can appreciate, if theta is 0.99, okay rate at which you will go to x star will be slower.

If theta is 0.1, 0.1 rise, k will go to 0 very, very fast. Iterations will converge very fast. So, what is the contraction constant, we will decide how fast you converge the solution, okay. So that is another message this theorem gives. The last message is very, very important. This is very, very important message, it says that I do not have to start from x not. See, we were talking about here.

**(Refer Slide Time: 51:57)**



So this is my x star and say this is my s not, okay. I am going to start my iterations only from x not, if I happen to start my iterations from some other x tilde not in the same ball, okay. As I said, you know she might take a different guess, he might take a different guess, he might take a different guess, okay. As long as those guesses lie within this ball, all those sequences will converge to the solution, very, very important, okay.

There is no unique initial guess. If you are in the region of convergence, any initial guess in that, if you give a good initial guess in that region, you are ensured to converge, okay. Sequence x tilde k generated by x tilde k+1 = G of x tilde k starting from any x not belonging to this region, okay. Where, I just continue this here, okay. So, if I were to start from any other initial guess than x not, okay.

As long as G is the contraction in this region, okay. I am guaranteed that the sequence will converge okay. All the concepts are important. Why Banach space? Any sequence that you start from any initial guess should remain within the space, very, very important, okay. Next thing is, we have this operator which maps this ball into itself. Then, it should be contraction, okay.

If it is a contraction, if all this conditions are met, these are sufficient conditions, if this sufficient conditions are met, you are guaranteed to get convergence to the solution, okay. So, this is the famous theorem called contraction mapping principle or contraction mapping theorem. There are many, many variances of this and I will just present to you one particular variant, which is easy to understand and very, very powerful.

We will just look at one or two examples briefly in the next lecture and then move onto the next topic. We cannot spend too much time on this because I will have to take many lectures if I really go into prove in this theorem, getting more insights, but what I want to do here by this one lecture is to just sensitize you that you know how do you look at the convergence properties of nonlinear algebraic equations, okay.

One simple message that you can carry is that look at the knob local Jacobian or G of x, okay. If that is not < 1, maybe you should try to make it < 1, so that you know you can ensure convergence and so on.