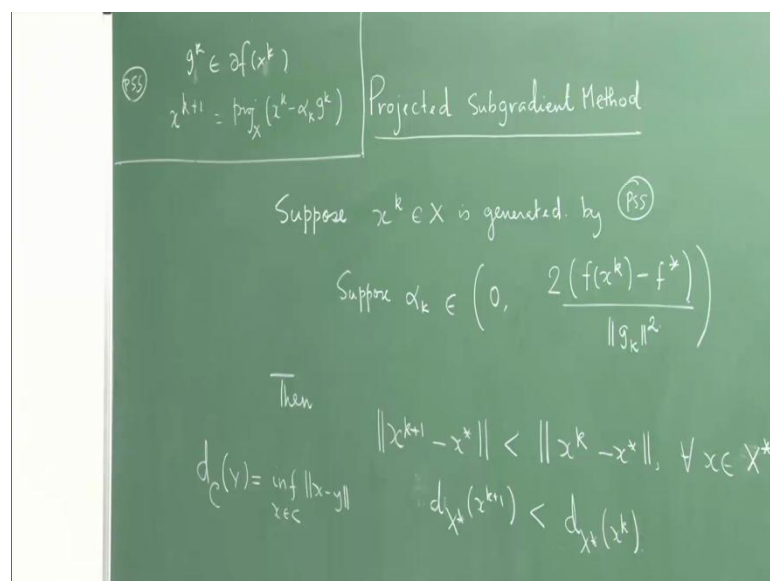


Foundation of Optimization
Prof. Dr. Joydeep Dutta
Department of Mathematics and Statistics
Indian Institute of Technology, Kanpur

Lecture - 35

When you begin a course like this, you have ambitious plans to tell the viewers a lot of things; you keep on doing it, and planning more and more adding things that you actually like yourself, because most what a person can do at the most is to, tell in many ways a story of his love, so you would really talk about what he actually likes in the subject. But suddenly one realize is that we sort of video courses are marked by certain amount of fixed lecturers, and then suddenly one realize is only at the end of the course, today is the last part 1 lecturer. So, it is a good idea that we would do some advance topics as we were discussing yesterday.

(Refer Slide Time: 01:28)



And we will talk about the projected sub gradient method; I am not going to write down all the details what I had written down yesterday, because it was already done. Now, tomorrow possibly I have to wind up this course with giving some very major ideas, but lot of things in optimization which I am not been able to speak here. For example, what happens if the data is a corrupted by noise that is a random variable which has coming to the data chance based input.

Around which the designer who ever in the special in the engineering stream has to know control on. What we can probably have some idea, from the experience you can have an idea about what sort of noise has coming what sort of probability chosen such random variable agrees to. So, the presence of noise of stochasticity is a very, very fundamental thing in optimization, such a thing cannot be a part of an elementary course, because one needs to be trained in much more things, in order to handle stochasticity in optimization.

But, for engineering problems stochasticity is a very, very big requirement, it is a fact that comes in naturally, we are not, we are not discussing this stochasticity here. So, the question is that these issues remember, that whenever I have a algorithm and I always trust that you need to show that, your algorithm works by showing the convergence analysis, does not mean that that algorithm is superbly efficient that that, once you start with the starting point it will rush towards the solution.

For example, the Stephen's (()) method that we discussed long back in the course, is a very, very good one and has a very good convergence for the quadratic, he is very good convergence for for example, a very good convergence for the convex cases. But, the problem is the following that it may not rush towards the solution it can be very slow, so and the so these sort of critics in that are convergence analysis, does not lead an algorithm to be very good enough for a given problem; and these criticisms actually are valid.

But, these convergence analysis also brings in to light, the quality of the algorithm itself, the nature of the behavior of the optima themselves, the nature of the behavior of the iterative iteration sequence; it tells you a lot about the nature of a particular algorithm, which has huge qualitative value. So, for example, I would like to show to you at least through this, so this is the original paper of written by Boris Polyak, called the minimization of minimization of un smooth functional.

So, it was published in US resort, and then later on it was you know translated to English by D Brown, but this paper and even now read the name, minimization of they did not write non smooth functional, but they have written un smooth functional; so that was a English translation. And it talks about it gives a unified method for both the smooth case and un smooth case where the Polyak step length that we had yesterday discussed, had actually been considered.

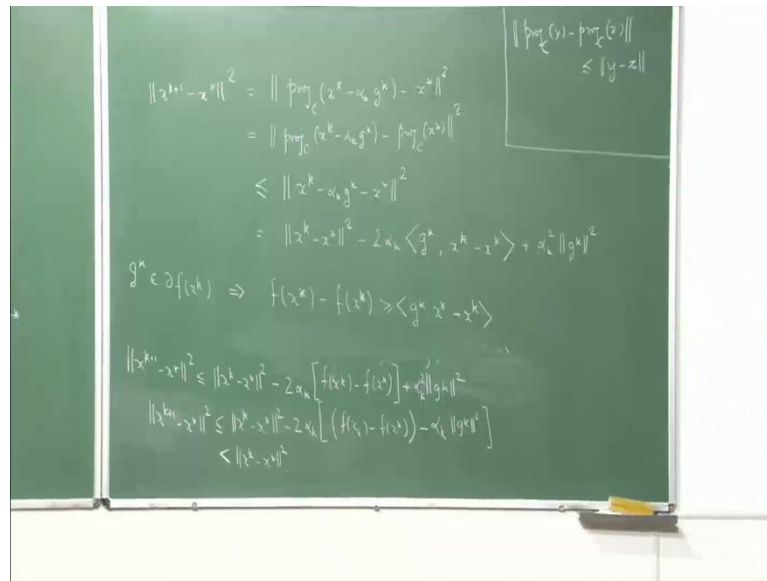
Now, let us suppose we have a sequence generated by the projected sub gradient method, the projected sub gradient method is, if you want to be more were you write, has to be written like this sorry, and x_{k+1} is projection on x . So, this is the projection projected sub gradient iteration scheme, and that is what we have going to use. Now, but suppose we had given that x_k , which are in x is generated by the projected sub gradient scheme.

Generated, if I call this as PSS scheme, Projected Sub Gradient scheme generated by PSS, so suppose x_k is generated by PSS. And your step length α_k , suppose now your step length is chosen like this, so it is a number lying in the open interval, so at every k your step length will change depending on your choice of g_k . Of course, here we are assuming that we somehow have come to know about the infimum value, but we do not know the x^* at which the infimum value would be achieved, and the whole process is to find that x^* . The whole process is to find that x^* using this scheme of iteration.

If this happens, then the most fundamental fact about the projected sub gradient method, that it does not tell you that the functional value is reducing, only tells you that the $k+1$ th iterate is much nearer to the solution x^* . So, whatever solution x^* you take, this fact is always true sorry, sorry capital X^* ; so, distance of x_{k+1} from the solution set is always little less than the distance of x_k from the solution set that is the mining of this.

Because, now you can take infimum over x^* on both sides to come to the conclusion, that distance a d of x_{k+1} distance function, this strictly less than the distance, so construct the infimum. The distance of the set see if I want to, to remind you, distance of a set over a point y outside a set C is given as the infimum of $\|x - y\|$, where infimum is taken over all x in the set C that is in that, so that is what happens.

(Refer Slide Time: 10:11)



So, I will prove this fact to you, so how do I write this one, so here is little bit of art of doing algorithmic analysis, so this estimate is very important; this actually brings out the character of the projected sub gradient method. So, this whole square is x^k plus 1 is nothing but, the projection, a very important fact about projection is, is it Lipschitz behavior. But, we have not spoken anything about Lipschitz functions in detail in this course.

So, we would say that if, so for any projection mapping, so you take the projection on the set C , then the norm of projection of C , so you take any any y and z in \mathbb{R}^n , this is always less than equal to y and z can be in C also. So, I can write this thing as, and where x^* is in C , the projection of x^* on the set C is going to be x^* itself. And that is what you have sorry (()) must be inequality, but the important fact is not just that there is an inequality of the distance, that given take any x^* this is always maintained, but when you take the infimum, you have the inequality here, less than equal to (()).

Now, once you have this fact, then you can apply that result to show that this is nothing but, now I have to take the whole square, because I have squared both sides. Once I have done this, then I can open the norm by writing this as I will block this two together, and have, so this is how you open the norms. Now, once you have opened the norms like this, the next step is the following, now observe that g^k is in $\partial f(x^k)$, so g^k is element of $\partial f(x^k)$.

$f(x_k)$ that would imply $f(x_k) - f(x^*)$, is x^* sorry, make a mistake $f(x^*) - x_k$ is g_k , $x^* - x_k$.

So, what do you have is minus, so if you take the negative, so what you would have if you just take this to this side you will have $g_k(x_k - x^*)$, take it to this side you will have $f(x^*) - f(x_k)$ sorry, so you will have $f(x_k) - f(x^*)$. Now, we are going to use the sub gradient inequality here, so I will just write down inequality again, so $x_k + 1 - x^*$ whole square is now less than $x_k - x^*$ whole square, you can take club in the minus here. Once you club in the minus here, it will become plus here, and $x^* - x_k$ and that will immediately give you that it is less than $f(x_k) - f(x^*)$, it will give you $f(x^*) - f(x_k)$.

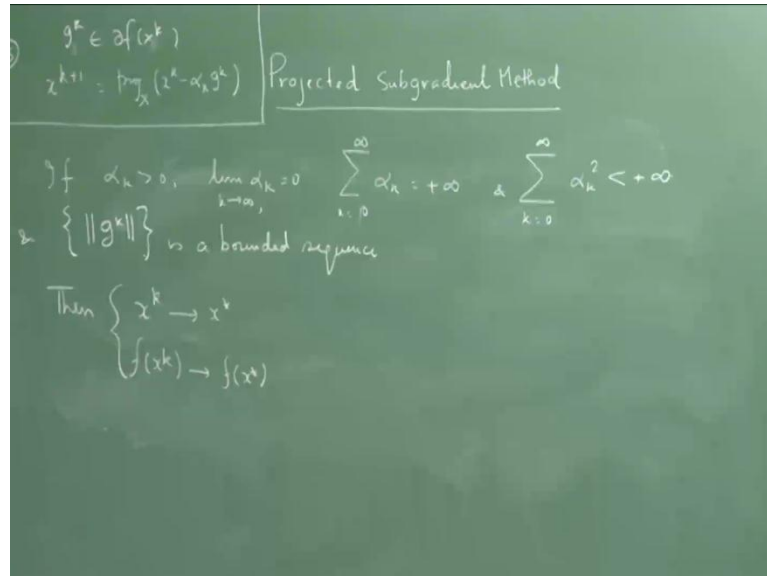
So, you (()) just bother about the last this one, this is the same thing, so this means I will have $-2\alpha_k(f(x_k) - f(x^*))$, which I can write as $f(x^*)$, (()) plus α_k^2 norm α_k this. But, this one $f(x_k) - f(x^*)$ is bigger, if x^* is smaller, so this is positive, so and this is how this when you and this is positive, so this whole thing is negative, so this is and now α_k , so this goes off; so this is less than this plus α_k^2 . There is no doubt that α_k is less than this quantity, so this is there is no doubt that α_k is less than the square of this is positive.

And now here what we have to do, is now we have to use the fact that α_k is less than this quantity, so you have $x_k + 1 - x^*$ whole square is less than equal to, now you have to take in to account this one. So, you write $x_k - x^*$ whole square less than, now you have to come with come up with the α_k , so maybe you can take negative twice α_k which is out; now you have $f(x_k) - f(x^*) - \alpha_k$ norm g_k square. I would tell you to finish this calculations by putting in the values of α_k , to show that this is norm less than norm $x_k - x^*$ square.

Now, once this is done we would like to go, and so this is also a motivation for the use of the Polyak step length, so if an α_k is chosen of this form, and you are the norm of sure α_k shown between these two numbers, then this is actually happening. In fact, this is strictly in fact this becomes strictly negative, so because α_k is not less than or equal to α_k is strictly less than this quantity, strictly less this, so this becomes strictly less. So, if an α_k is taken to be equal to that this become less than or equal

to, so it make sense to use that particular choice of alpha k, and that is from where this Polyak step length is important.

(Refer Slide Time: 19:50)



Now, we are going to try to prove or give a sketch of the proof of what we intended to show yesterday or we mentioned yesterday that, if the assumption one means if the step length, if alpha k is strictly greater than 0 limit of alpha k is equal to 0 as k tends to infinity. Summation alpha k k is equal to 1 to infinity is plus infinity, and summation all of this has to be satisfied, the step length has to be chosen k equal to 0, this k 1 by k plus 1.

So, if this happens, if g k is bounded now, because I have not given you the mathematics of the fact that the sub differentials are, if you take finite valued function from one and two are you, the local what they something called, local boundedness of the sub differential; and that would lead to the fact that norm g k is bounded right. So, g k is element of g x k, and x k is, no we cannot say that sorry I made a mistake in the last class last class also, let me admit that mistake.

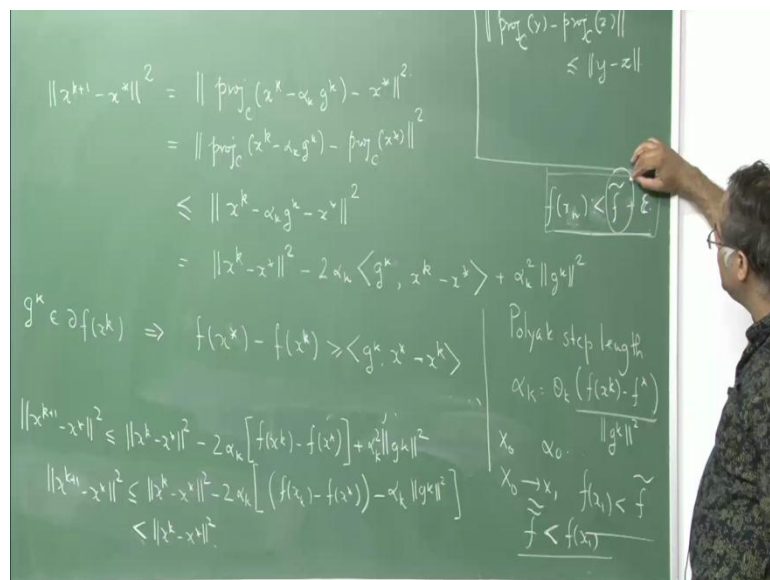
We have not yet proved that the x k is going to x star, if x k goes towards x star, then we can say that g k is bounded, so that is what we really have to prove that x k is going towards x star. So, g k has to, the norm of g k this is assumption, it is a bounded sequence then we are proving that, actually you have unless you have x k going x star you cannot use that, whatever I am calling as local boundedness property. So, please

now, what you about the local boundedness property at all, but remember when Polyak wrote those things these are essentially done for functions which are extended valued of course, x^k are all in the domain.

Because, they are in a feasible set, when over which has find it, so must the domains intersection with the feasible set must be known empty, so you need not get too much poked away with it, but we have to take this thing, but if you take the Polyak sequence you do not have to bother about this boundedness of this norm g^k . If this is assumed then x^k goes to x^* and obviously by continuity of a function f , if it is of course, if you do not have continuity have to show that, but if you have continuity of the function f then it is automatic this goes to this, so these two things happen.

But, our case this will happen will imply that this will imply this, so you forget about this word local boundedness that I have used in that lecture; so this is assumption I am making. So, you might be asking me why I am proving this result, and not the result with Polyak step length, see Polyak step length has this issue.

(Refer Slide Time: 23:15)



So, what is for the Polyak step length, there is a fundamental assumption when the Polyak step length is something like this which, so f^* star the infimum value sorry, it becomes (()) step it is changing, because of this. And we know there are particular choices of θ has to be done that we have written yesterday, now the problem is the

fact that I know f^* , the in real computations you cannot have f^* , you do not know f^* of course, because you are supposed to find f^* .

Because, you do not know f^* , you really it cannot be inverse problem that I give you the f , I give you the f^* then tell you to find x^* , then you use the Polyak step length straight away. But, if you do not have such a thing and what you can do that you can take an arbitrary estimate of f^* , and then try to say from x_0 go to x_1 . If $f(x_1)$ is strictly less than f^* then of course, your function values that f^* is not a good estimate. So, what you have to come in to a situation, so you what you do, so you reduce the value of f^* , you reduce it below $f(x_1)$.

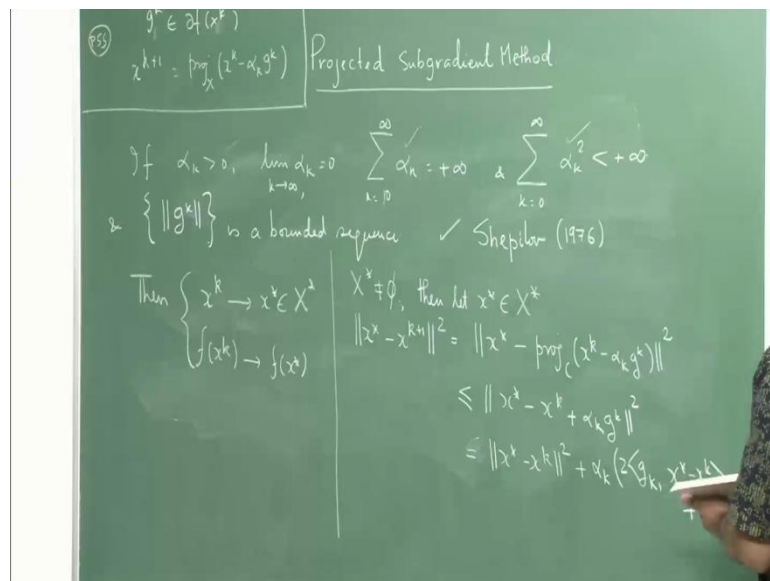
And then you again try, so you try for some times and see that you have got $f(x)$ value less than equal to less than f^* , so you can try with an x you take that x_1 and then try to the new estimate, you can try from going from x_1 to x_2 with the new estimate of f^* , and the step length. So, from here I start with the point x_0 , and then I calculate α_0 and then I go from x_0 to x_1 by this standard projection. So, now I check whether if $f(x_1)$ is strictly less than my estimate of f^* , then possibly this not a good estimate of f^* is, f^* cannot be strictly less than $f(x_1)$.

So, then I have to choose something which is say f which is strictly less than $f(x_1)$, and then come from x_0 to, and then take that x_1 as a starting point, and then I try to get to $f(x_n)$ or may get to $f(x_2)$. So, what I essentially would require at the end is that ok, I must have my $f(x_n)$'s of this form, it would be enough to have this situation in computation for some given ϵ greater than 0. This is this is the sorry yeah, if f is at infimum then if I do this and x_n at least should be of this form, so for some n , x_n you should immediate ask given fix ϵ you can have it in this form, so f^* is, f^* at is the infimum.

I made a mistake $f(x_n)$ is sorry, I made a mistake it is infimum, it should be less than or equal to rather strictly less than $f^* + \epsilon$. So, if this happens then we can take this f^* , this is the definition of ϵ , if definition of infimum, if f^* in some sort of an estimate of the infimum suppose this is the infimum, then for some ϵ greater than 0 I can find an n , so that $f(x_n)$ must be this. So, now if I when f^* and some f^* I can I have an x_n which is this, and that and if I lock my $f(x_n)$ value is decreasing, then this f^* would work.

So, this certainly a complicated thing, and so here Polyak step length is largely a tool which would give me theoretical results, but the problem here again is that these step lengths. And immediately you can find it a step length satisfying this, but the problem again would be this, who check the boundedness of the norm of g_k , how do you know that they are actually lying between something, so these are questions that coming. So, everything has an advantage everything has a disadvantage, because you know once you come from to sub gradient issues, but this can also be used in the gradient thing.

(Refer Slide Time: 28:43)



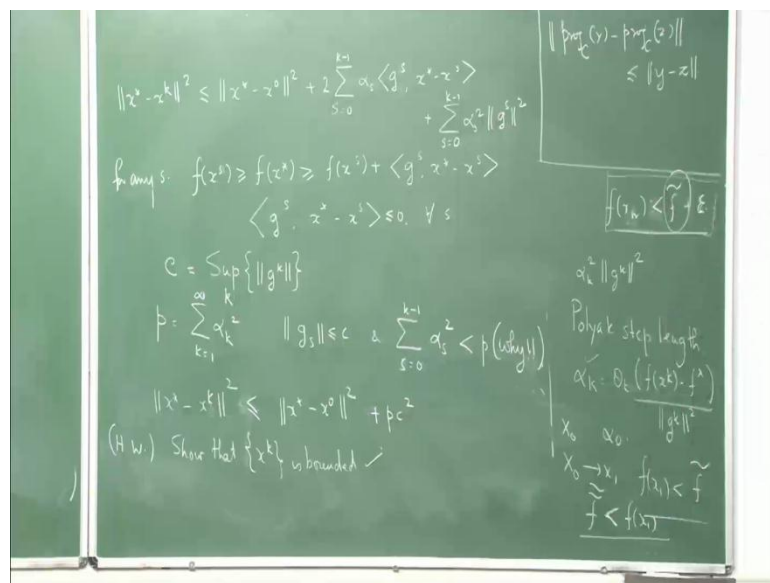
So, but only advantage in the Polyak case is that you do not have to assume the norm g_k is bounded that that is the plus point, any way let us this was the proof given in 1965. And let me just tell you who was the author, this proof was by Shepilov 65, the proof that we are going to give is by the Russian mathematician Shepilov, Shepilov gave this proof in 1976. So, Polyak result is in 1969, so Shepilov improved it in the sense that is easy to find it, but this is again now draw back. But, luckily if this is happening, so this also happens, so practical purpose is this is pretty, the Shepilov result is pretty pretty interesting in that sense.

So, we try to do some proof of this, some outline of the proof I would not say the full proof, but I will just go and do some outline, give some idea us to how the proof can be achieved. But, that would constitute what would one call as a art of convergence analysis, so again you start with this fact that you assume that x^* is a solution, so x

star, so this goes to, so we have assumed that x^* is non empty. So, we have assumed that x^* is not equal to ϕ , then let x^* is element of x^* ; now compute this distance, from x^* compute the distance of x^{k+1} square of the distance rather that is again the same as what we have done earlier.

So, this again by that same property written there, at that end of the board this property, would lead to the fact that minus minus plus this minus and this minus will be plus, and so this is what you have. And then if you write down open the square, this is exactly equal to, so α_k in to, I am writing this I am not doing too much of stepping jumping a step. So, basically I have $2\alpha_k g_k$ in a product $x^* - x^k$ and of course, α_k^2 plus norm g_k square, so we are not combining those two terms and writing it immediately to get.

(Refer Slide Time: 32:33)



If I use this thing repeatedly what I have is the following, I have that $x^* - x^k$ whole square, so I I I repeat I start from x^k go to x^{k-1} , so here it will be x^k then x^{k-1} , so and so forth. So, and then keep on repeating to get the starting point, so obviously will be up to $k-1$, this things repeat at applications of the cone. So, s is a (()) my variable instead of not g_s here, what you have, now f of x_s for any s for any s guaranteed.

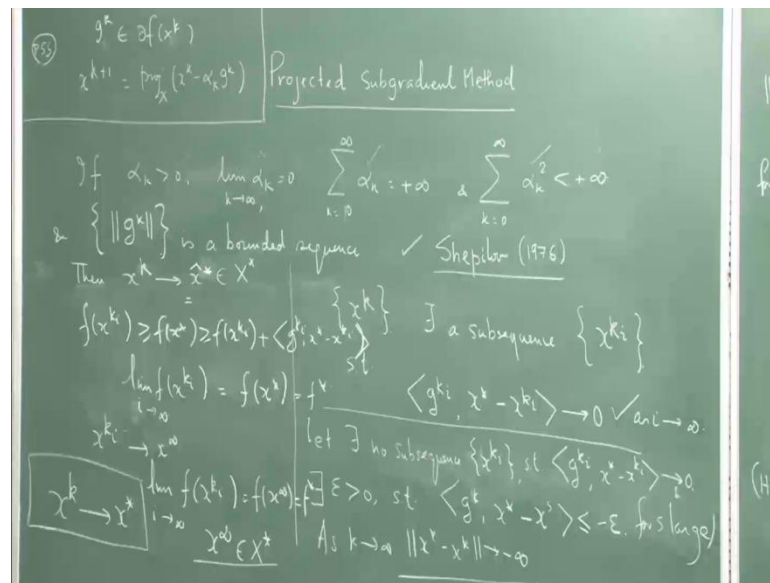
Now, define c , because you know that this is bounded this norm g_k is bounded then define supremum of norm g_k over k I define this, now say p I shift that the infinite sum,

because this is strictly less than infinity, this is p . And this is obviously less than c , and I would have sorry, why this happens I would this is a very basic question in analysis, so I would just ask you to figure out this, I would not repeat why such a things happens, this is the very basic question in very basic basic calculus basic analysis. So, because we had the advance level at the end of the course, certain advance level, so this is you can figure out why.

Now, then going back to what I have here, I can write that $x^* - x_k$ whole square is less than or equal to right, this part is less than equal to 0, so this negative this is goes I have this part, and this part. So, it will give me sorry, this is x_k , this one is this plus, now this can be written as strictly less than, because this is strictly less than this this, so there will be a norm is less than equal to c . So, product is strictly less than unless of course, there is a 0 fact (()), so this will simply show you that from here home work show that x_k is bounded.

So, this is a very very important step in convergence analysis to show that the sequence that you are generating the is bounded. Because, then you are guaranteed there is a convergence sub sequence, so the idea always is that you show that there is a convergence sub sequence which will have a limit point, so that convergence sub sequence has a limit point. And then show actually the whole sequence goes to that limit, so and the whole sequence is actually going to that limit, and hence and show that the limit is having some properties like this that is element of the solution set. So, this can be showed as a home work, this is too easy for us to unnecessary detail on, because you have to understand, because x not is known, $x^* - x$ not is actually a fixed quantity, and this is a fixed quantity; the crucial part of the proof lies here, which we will now analyze.

(Refer Slide Time: 39:14)



Crucial part of the proof says the following, but Shepilov shows is the following that in this sequence x^k , there exists a sequence, there exists a sub sequence x^{k_i} such that that now what you are going to prove that there exists, in x^k there is a sub sequence x^{k_i} , there exists some sub sequence or sub sequence such that g^{k_i} this goes to 0, the proof this though it is a very, very crucial step in the proof of Shepilov theorem, Shepilov analysis; this is a very, very interesting thing in the sense that its proof is not very difficult, because we will use one of the most important tool of the mathematician tool works proof by contradiction, we will say that let us now assume that, there is no such sub sequence in x^k for which this will happen. Let there exists no sub sequence such that, sub sequence x^{k_i} , but you know that this result is true, so though is less than equal to 0, the sub sequence value cannot tend to 0, for whatever be the sub sequence.

So, given any epsilon greater than 0, so the not given any epsilon there must exist some epsilon, so there must be a negative number which, which will bound this value, so this value cannot go above that negative number and move towards you. So, which means there exists epsilon greater than 0 such that, $g^k \cdot x^* - x^k$ must be strictly less than minus epsilon, this is this is what you have. At least for sufficiently large values of x this must be true, that is for some finite number of values it may go beyond that, but after that for s for s large for s large sufficiently large it means in this has to be the case otherwise, it can go to 0.

But we are telling that there is no such sub sequence for which it goes, so for s sufficiently large or this s that you are generating or the, so whatever sub sequence you take. So, for sufficiently large index of x_s in the in the sequence, because we are choosing x_s from the same sequence x_k , this has to be true, now if this has to be true what is the consequence. So, remember this is a very, very crucial thing, so we will just write here that will just write here that this statement of a of the result final result sorry very bad mistake.

Now, if this is true I just want to remind you the sequence x^* that I that x_k I generate need not go to that x^* that I have chosen here, so I will put here that this as some x^* x^* x^* , x goes to x^* , so it goes to some x^* , x^* . So, I started with some x^* in x^* , because x^* is non empty, but that does not mean that the sequence I generate goes to this one. Now, what would be the consequence if this is happening, so from here the consequence is this is less than minus epsilon, so why I say that this is a very crucial step, because now you will see the use of this results, the use of this we have used this results, but we have not used this fact, now because you have this, we have essentially used this and this, but we have not used this result, this is not used.

So, this has been used, this has been used of course, α_k is any where greater than 0, this has not been used, we will not show the use of this we will just keep it, we will not bother about at this moment, this has not been used. So, if I use this, if I say that this is less than equal to some minus epsilon, then this as k becomes large these goes to infinity, so this whole thing goes towards minus infinity as k tends to k tends to infinity. So, as k tends to infinity this would imply that $x^* - x_k$ goes to minus infinity, which is obviously not true.

Because I have said that we have already showed that this is this is a bounded thing, so I cannot so this will contradict the fact that this is going, so what I have said is this fact that there exists no sequence is not true, and there exists a sequence where this takes place. Now, once this is done, when you know this fact, then if you know that there is a sub sequence, now you come here you put x_{k_i} here, x_{k_i} here, and this x_{k_i} and x_{k_i} here, and you take the limit, so from here what I will now do, I will take this inequality and in this inequality.

Let me just have a use this fact here, so I am taking now this inequality, so I am taking this inequality and what I am proving, I am proving that I am proving this following fact that $f(x_k)$ is greater than $f(x^*)$ is greater than $f(x_k) + g(x^*) - x_k$. Now, as x_k sorry, $g(x_k) - g(x^*)$ as k goes to ∞ , as k as i goes to infinity. So, as i goes to infinity which means what, that i goes to infinity this is going to 0, which means the limit of this by using the sandwich theorem, limit of $f(x_k)$.

Now, which is which I can call as f^* , now x_k is a bounded sequence, the bounded sequence $f(x_k)$, there is a convergence sub sequence, we assume without loss of generality it is at x_k goes to some x_∞ . So, then I can immediately write by the continuity continuity of the functions f that, limit of x_k , i tending to infinity is $f(x_\infty)$ is equal to f^* . And x_∞ , is obviously an element of x^* , because this value is equal to f^* the solution, the optimal value; so x_∞ is now in x^* .

And my job the problem of whole proof would end, once I can show that x_k goes to x^* sorry x_k , the whole sequence x_k will go to x_∞ . So, this is x_∞ is my x^* , so basically now I have the proof will end, once I prove that x_k goes to x^* , and this part of the proof I will not prove, because we are running out of time. But, this is a whole glimpse, what I have done here is to tell you that o_k this is the way things are done, so there is a crucial point where you start applying this idea. And this idea, this idea would be used, when you are trying to prove this, you see that we are using everything, and that that is very, very important.

So, I have just told you the crucial thing, so we have shown that, we have we have generated a solution, we have generated a sub sequence of x_k from which I have some behavior from, which goes to some solution. And actually we will show that, no not x^* sorry I made a mistake, the whole thing actually x_k goes to x_∞ , please take care of this. So, that is what we have to prove, at the end which have not proved and with this we end our last today's lecturer, which is the last one, basically for the course and tomorrow we will sum up what we have studied.