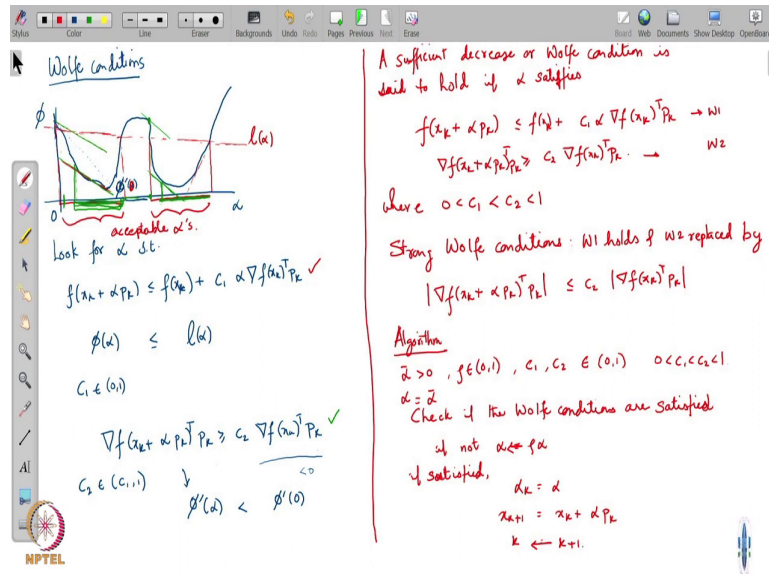


Optimization from Fundamentals
Prof. Ankur Kulkarni
Department of Systems and Control Engineering
Indian Institute of Technology, Bombay

Lecture - 19C
Line search algorithm and convergence

(Refer Slide Time: 00:16)



A practical line search type algorithm the algorithm stands to work like this. So, let me write this out. So, you start usually with some say an alpha bar greater than 0. You choose a rho between 0 and 1 and I will tell you why this rho is needed. You choose your C 1, you choose your, C 1 and C 2 and then what we what one does is you is that you search overall.

So, what one does is you try you we can say check if the Wolfe conditions are satisfied. So, at each iteration check if the Wolfe conditions are satisfied. Check if the Wolfe conditions are

satisfied. If not, you try an alpha. So, you in right initialize alpha equal to some alpha bar, if not, change alpha to rho times alpha.

So, what you are doing is you start with a large enough alpha and then you keep back tracking to see where you can to see what would be if your Wolfe conditions are satisfied right. You and you can keep repeating this until eventually the Wolfe conditions are satisfied. If satisfied you would define alpha K as alpha define alpha K as alpha and the next iterate x K plus one is defined as x K plus alpha P K, alright. And then and you take K as is we set K to be K plus 1.

(Refer Slide Time: 03:44)

$x_{k+1} = x_k + \alpha_k p_k$
 where p_k is a sequence of directions
 if α_k satisfies the Wolfe conditions.

Convergence of line search algorithms.

$$\cos \theta_k = \frac{-\nabla f(x_k)^T p_k}{\|\nabla f(x_k)\| \|p_k\|}$$

Thm: Consider any iteration of the form $x_{k+1} = x_k + \alpha_k p_k$ where α_k satisfies the Wolfe conditions w_1, w_2 . Suppose f is bounded below ($\inf f > -\infty$) f is continuously differentiable.

Suppose the gradient ∇f is Lipschitz continuous on an open set containing $L = \{x \mid f(x) \leq f(x_0)\}$ where x_0 is the starting/initial iterate, i.e. $\exists L' > 0$ s.t.

$$\|\nabla f(x) - \nabla f(\bar{x})\| \leq L \|x - \bar{x}\| \quad \forall x, \bar{x} \in L.$$

Then
$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(x_k)\|^2 < \infty$$

In particular $\lim_{k \rightarrow \infty} \cos \theta_k \|\nabla f(x_k)\| = 0$

If p_k is chosen so that $\cos \theta_k = \epsilon > 0$
 $x_k \rightarrow x^*$ then $\nabla f(x^*) = 0$.

Simplest approach $p_k = -\nabla f(x_k)$
 Steepest descent algorithm

So, what happens then is so, your algorithm then effectively then takes you through these iterations, where x K plus 1 equals x K plus alpha P K, where P K is a sequence of directions and alpha K satisfies the Wolfe conditions at each K. So, this is basically so, what we get therefore, is an is a sequence of iterate such that that are given by this recursion x K plus 1

equals x_K plus $\alpha_K p_K$, where p_K is a sequence of directions and α_K is no satisfies the Wolfe conditions.

Now, what I have not told you yet is what how is this p_K to be chosen. So, we have been we will been silent on this particular topic about how we how about what the choice of p_K is. So, the p_K is to be is can be chosen in a variety of ways and I will give you a general result here, but the main thing is that the p_K should be a direction in which the function decreases.

So, the simplest thing you could do for this sort of purpose is to simply look for a $d p_K$ to be the negative of the gradient of the function. So, locally we are assured that whenever the in that if we take a step in that direction if we move infinitesimally in that particular direction the function is guaranteed to decrease.

So, p_K can be taken to be the negative of the gradient of gradient of the function, but it that is not that it is that is not the only choice, one could do many other things to choose the choose the p_K . In fact, variants of this give you gives rise to various different types of algorithms itself ok.

So, I will now also mention to you is a sort of a generic sweeping result which ensures which tells us what sort what sort of p_K 's to line search algorithms actually to line search actually a line search algorithm are actually converge ok. So, now, define. So, that brings us to this topic of convergence of line search algorithms.

So, define \cos of θ_K as basically the cosine between the negative of the gradient of the function and p_K . So obviously, this kind of quantity is because so, its I am taking the inner product in the numerator negative of the inner product and dividing by the norms of these two vectors.

So obviously, this is well defined only when the when these two vectors are not 0. So, p_K is obviously, a direction we are choosing. So, it is a non-zero direction and gradient of the

function is if it should be non zero, then this quantity is well defined ok. So, the so, what is the; what is θ_K here?

θ_K is capturing the angle between P_K and the steepest possible descent direction you could pick. The direction that you could pick which is which gives you the steepest decrease steepest most decrease in the vicinity of the function, right.

So, this is the so, θ_K is capturing that the angle between these two. So, the theorem which we will not prove, but I will just mention to you is this. So, consider any iteration of the form; consider any iteration of the form $x_{K+1} = x_K + \alpha_K P_K$.

Now, where α_K satisfies the Wolfe conditions, W_1 comma W_2 . If now suppose f is bounded below. Bounded below means, it has the infimum of f over \mathbb{R}^n is greater than minus infinity and is continuously differentiable. Suppose, the gradient $\text{grad } f$ is Lipschitz continuous on an open set containing the set containing the set L .

Let us call this set. What is this set? This set is called the level set its the set of axis for which $f(x)$ takes value less than equal to $f(x_0)$, where f where x_0 is your starting iteration; starting iterate starting or initial iterate. So, now what does it mean? So, the suppose the gradient f is Lipschitz continuous on an open set containing this level set, where x_0 is the starting iterate. What does this mean?

That is there exists an L L dash say greater than 0, such that gradient of f at x minus gradient of f at \bar{x} open set let us call this open set N less than equal to L times norm of x minus \bar{x} for all x, \bar{x} in this open set N ok. So, we have we can take we can take P takes P_K to be any set of directions like and choose α_K such that the Wolfe conditions are satisfied ok. And suppose the function is bounded and continuously differentiable and the gradient is Lipschitz in an in the set that we are considering.

So, what is this set L ? L is the level set means that it is from starting from take all the x 's that give you a better value than the one that you are starting with this the one you are starting with is x_0 all the x 's that give you a better value a lower value than what you are starting

with that that entire region is called a level set that is your set L ok. And so, we want we are assuming that the gradient is Lipschitz on that set.

Then what does it say? Well then it says then $\cos^2 \theta_K$ times norm of gradient $\times K$ square this summation is finite. This whole sum is finite. Now, what does this mean? So, the claim is that if you can if you choose your iterations this way if your function has these properties and you choose your iteration this way in the way that is indicated and your α_K satisfy the Wolfe condition, then this particular condition this has to be finite.

Now, this looks like a like a bizarre technical conclusion, but actually it it is says a lot in one in one sentence. So, since this summation this is an infinite sum right. So, if this infinite sum is finite what is that? If that infinite sum is finite then it means that the limiting value of limiting term here should be going to 0.

So, in particular in particular limit as K tends to infinity $\cos^2 \theta_K$ times norm square of this is equal to 0. So, it says that well your, the this limit goes to 0. Now, if you look at this limit what is this limit? It takes it has two terms here. It has the norm of the gradient of the function and it has the \cos^2 of the cosine of the angle between the negative gradient and the direction P_K that you have chosen.

So, now this tells you something quite nice and powerful. It tells you that if I can choose my, if I choose my P_K in such a way so, if P_K , so, what does this mean? If P_K is chosen is chosen so that say $\cos \theta_K$ is say always equal to some epsilon ok say minus epsilon ok, rather since you have already a minus sign there this is equal to epsilon say. So, if I choose my P_K so that the cost θ_K is always equal to some constant epsilon, which is positive.

Then in that case in this limit here; in this limit here this $\cos^2 \theta_K$ would always be epsilon and it could it can jump out of the limit. And in that in the and then what we are left with is just is that the gradient of the limiting value of the gradient of the function should be 0, which means that x_K converges to if x_K converges to some x^* then the limiting value of the gradient of x^* should be equal to 0 right.

So, if $\mathbf{p}^T \mathbf{k}$ is chosen in such a way that you are making an acute angle with the negative of the gradient, so, it has some component along the negative; it has a non zero component along the negative gradient.

Then you are guaranteed that the sequence of iterates actually get you to a point where the gradient becomes equal to 0; that means, you are satisfying the necessary conditions of optimality. Now, this means that. So, this it is important here that this angle is this the cosine here is some epsilon equal which is positive and remains positive and so, the way we have done this is by choosing an epsilon that is independent of K .

If epsilon also depends on K and starts decreasing to 0 then in that case the limit of this product being 0 does not let me conclude that the gradient is equal to 0. It could well be that you know this product the gradient still remains positive and yet your, this limit is going to 0.

So, the so this condition is effectively telling you that the way to ensure that your gradient vanishes is by making sure that the, this cosine remains bounded away from 0. So, you; that means, you should continue to make if a non vanishing angle the that $\mathbf{p}^T \mathbf{k}$ should continue to make a non vanishing angle with the negative gradient whatever that gradient may be ok. So, $\mathbf{p}^T \mathbf{k}$ if it continues to make this non vanishing angle with the negative gradient then the gradient itself will go to 0, right.

So, the simplest way of ensuring that is the simplest way is to take approach to doing that is to take $\mathbf{p}^T \mathbf{k}$ to be the negative gradient itself. So, then you are collinear with the negative gradient and in that case you would obviously, make your yeah. So, in that case this would actually go down to the angle would go the angle will the epsilon in that case will always will be 1 and then the gradient would be equal to 0. This kind of algorithm is what is called the steepest descent; steepest descent algorithm.

Now, steepest descent only ensures that you are taking the steepest the direction of the steepest descent at each point that may or may not be right for you in the long run. In the sense that the steepest what looks like the steepest descent at a particular point may not give you a sustainable decrease all the way down when you go further.

So, you have to take into account also how the directions of steepest descent themselves change. So, you so, the ideal way is to actually take into account also the curvature of the function and that those kind of that gives you a much richer class of algorithms that, but they are once again another class in this kind that I have mentioned.

So, they so, long as the, you when you are designing these sort of algorithms or when you are coming up with your iterate make. So, long as you make sure that your cosine of the angle is continues to be bounded away from 0 you should be fine and you your, you will be going to your iterates will take you to the grade to the to a point, where the gradient vanishes alright.

So, I will so, with that I will stop this lecture and we will continue next time.