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

Lecture - 21A L1 and L2 Penalty methods

Welcome everyone. So, we were in the previous lecture we were talking about Penalty Methods. Penalty methods are a way of converting a constrained optimization problem into an unconstrained one by penalizing the constraint. So, we introduce this additional function called the penalty function.

The penalty function was supposed to have the was required to have the property that it would be continuous it would be non-negative and it would be 0 entirely on the feasible region of the problem. So, it would be 0 everywhere on the feasible region of the problem and outside the feasible region it would be strictly positive right.

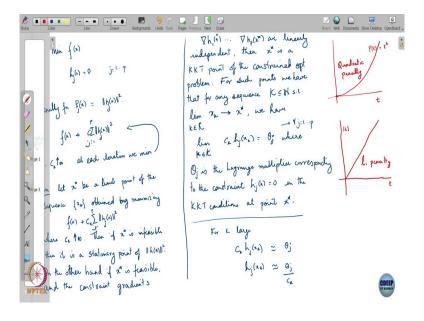
So, and using this penalty method penalty function we recast a constrained optimization problem as an unconstrained optimization problem; by removing the constraints and adding to the objective a penalty parameter times the penalty function. So, the new objective was the original objective plus a penalty parameter times the penalty function.

Now, one of the things we observe we found there was that as the limit if you computed the and the sequence of unconstrained minimizations by letting the penalty parameter go to infinity one of the things we found there was that this sequence, this sequence if it converges; that means, any limit point of this sequence converged to the global minimum of the constrained optimization problem.

Now, this is the extremely powerful technique because it allows us a way of essentially reducing any unconstrained optimize any constrained optimization to simply an unconstrained one. However, there are a few little drawbacks which I will just discuss.

One of the things that we find in a in the penalty in the in this sort of method is that the that one converges rather slowly to a feasible point meaning that the feasible in order to become strict actually feasible your penalty parameter needs to be extremely large right. So, let me let us see one kind of penalty method in which we are using what is called a norm penalty or a quadratic penalty ok.

(Refer Slide Time: 02:56)



So, consider this optimization problem. Suppose we are minimizing suppose we are minimizing this function f(x) subject to the constraint h(x) of h(x) equal to h(x). Now, one possible choice for the penalty function for this sort of problem since this is a equality constraint one possible choice for the penalty function is to simply penalize h(x) by the penalty function h(x) of h(x) can be taken as simply norm of h(x) of h(x) the whole square ok.

And then we can consider the objective which is f x plus c times norm h of x whole square. And the idea is that, we let us c increase to infinity increase to infinity and then we at and at each iteration we minimize this particular problem right.

Now, the where is the catch? Here, so, the catch is the following. So, let me mention this theorem. So, the theorem is the following. So, if so let x star be a limit point let x star be a limit point of this sequence x k. So, let me say for example, let us take c as c equal to c k and c k goes to infinity right.

Let x star be the limit point of the sequence x k obtained by minimizing f of x plus c k norm h of x the whole squared alright and obtained by minimizing f of x plus c k norm h of x the whole squared, where c k is tending to infinity is increasing to infinity alright ok.

Now, then if x star is infeasible then it is a stationary point of norm of h of x squared. So, if it is infeasible it actually ends up as a stationary point of norm of h of x squared. So, you end up actually minimizing not f, but or not f plus this, but rather eventually ending up minimizing this.

So, you end up at a stationary point of this ok. On the other hand, if x star is feasible and the gradients of the constraints. So, here by I did not mention this, but this is these are basically this could be even a vector of constraints. So, we will allow for this.

So, let us say let us write this as a j equal to 0, j equals 1 to P say and then we could write this as norm of this ok. So, this would be the penalty function for the kth constraint and this would be then a summation from 1 to p ok alright. So, I will just adjust this to allow for a vector and the constraint gradients. The constraint gradients are these are linearly independent, then x star is a KKT point of the constrained problem of the constrained optimization problem.

For any, for such points we have that for any sequence K N greater than such that the limit as k in K x k tends to x star, we have limit k in K c k times h j of x k equals we have that the

limit as k of k as k runs over the sequence capital K the limit of c k times h j of x k equals theta j where theta j is the Lagrange multiplier corresponding to the constraint h j of x equal to 0 in the KKT conditions at point x star.

So, what does this mean what the statement effectively is saying is that if you take the sequence of x ks that are obtained by solving the penalized problem and you let this you look at any limit point x star of this sequence then you have two possibilities.

One is that the limit point is actually infeasible for the original problem in which case it turns out that the limit point is actually a stationary point it is it in which case it turns out that the limit point is actually a stationary point of this of your penalty function. So, you end up actually at a stationary point of the penalty function and that has may have no in relation whatsoever to the to any solution of the original problem.

On the other hand, if you are feasible then it turns out that if your constraint gradients are linearly independent, say for instance if L I C K holds and then you actually end up at a KKT point of the constrained optimization problem moreover with this penalty function moreover you have you have this additional property that c k times h j of x k tends to approaches the Lagrange multiplier c k times h j of x k approaches the Lagrange multiplier.

Now, this is a particular property of this quadratic penalty that we are that we are considering. So, it this limiting value becomes the limiting becomes the value of the Lagrange ends up at the optimal value of the Lagrange multiplier at for which for the KKT conditions at point and at point x star alright. So, this is and this is true for every each equality constraint ok. So, this is for all this is solves for all P for all j in going from 1 to P all the P equality constraint.

Now, what does this mean? This means that effectively for k large this basically saying that for k large your c k times h j of x k is not becoming 0, but rather become coming close to a constant equal to the Lagrange multiplier right. So, which means that if I will just think of h j of x k itself h j of x k starts approaching a constant divided by c k.

Now, this particular this is this particular thing is a somewhat undesirable because what it effectively means is that h k is never actually going to become exactly equal to 0, which is what you would need h j h j is would has to be exactly equal to 0 for you to have feasibility.

So, you are never actually going to have a h j exactly equal to 0 unless c k itself becomes infinitive right. So, unless c k becomes infinite this is not going to exact work out exactly right. So, which means that you really need c k to blow up to infinity and the and. So, for if you terminate the algorithm at any finite iteration although in the limit you would end up at a solution.

But, if you terminate the algorithm at any finite iteration this you may not actually be feasible. In fact, you would in general not be feasible you would be feasibility will be of by a certain by a certain amount. Now, there are two ways of remedying this and let us I will talk about one particular approach first and then we will go to another approach. So, one particular approach is to change the penalty function itself.

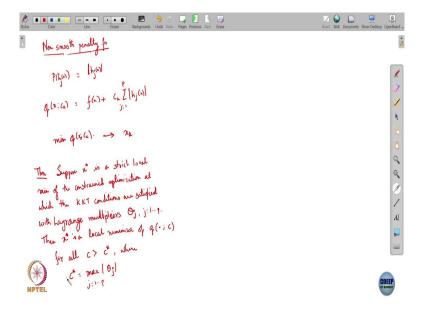
So, the reason this is happening is because your penalty function right now is actually quite smooth the penalty function that we are using has increases the quadratic penalty function tends to increase gradually; if you see the this this quadratic penalty function effectively if you see how it looks essentially looks like a gradual increase towards infinity that is the; that is the behavior of this function right t and P of t, so the penalty function P of t which is equal to t square.

So, it is not quite doing a sharp penalization near the feasible region. So, as a result your method tends to sort of become a little bit ambivalent about points that are closed that are around here that are close to the feasible region, but are but may be feasible or infeasible. So, it is not the penalty out here is not sharp enough. Now, the way one way of therefore dealing with this is to put in that sort of sharp penalty and that is what I will talk about right now.

So, the sharp penalty would effectively mean that you are putting you need to put an a penalty that goes from the from inside the feasible region to outside the feasible region in a very in a

dramatic sort of way. So, that if usually means that you are you would lose differentiability of the penalty function you can have continuity, but the function will not be differentiable any more just like the quadratic function smoothly increase you will not you are not going to you cannot expect that sort of behavior.

(Refer Slide Time: 18:30)



So, you are then looking for a penalty function that is not smooth right. So, that is what is called a non smooth penalty function. So, a non smooth; a non smooth penalty function. So, one a simple example of a non smooth penalty function is the modulus function so the absolute value function right. So, if you have so, you take the absolute value.

So, P of h of x as simply or h j of x as simply the absolute value of h of x instead of the square of this quantity. Now, this function unlike the one before the this is your quadratic

function quadratic penalty the one with where you are using the modulus or the absolute value that function would look like this.

Now, if absolute value being less would eventually will eventually this being linearity will eventually the not penalize as severely as the quadratic because the quadratic would penalizes in a much more dramatic way for larger values.

But for smaller values that the penalty in a in the absolute value or the 1 1 penalty this is also called the 1 1 penalty. So, the for smaller values of t the penalty here is going to be larger because for smaller values mod x was actually be smaller values of t mod mod t will actually be our absolute value of t would be less than t square right. So, this is the thing that we would end up we want to end up exploiting.

So, then what happen what we are then looking for is then we are the problem we are solving then is this problem q of x given c k which we are in we are going to define this q of x k minus c k as f of x plus c k times summation absolute value of h j of x this j ranges from 1 to P and the and at each step you minimize you minimize this q of this in order to get x k.

Now, the theorem that we can get from this is following. Suppose x star is a strict local minimum of the constrained optimization problem of the constrained optimization at which the KKT conditions are satisfied with Lagrange multipliers with Lagrange multipliers theta j for j in 1 j going from 1 to P. Then x star is a local minimizer of the x star is a local minimizer of q and c for all c greater than c star, where; c star is equal to the maximum of the Lagrange multipliers maximum of these Lagrange multipliers ok.

So, x star turns out to be a is a local minimizer is a local minimizer of this for all of the or for all c greater than c star. What does this mean? That if you found if you can find if you can set your penalty parameter to be larger than c star where c star is simply the largest of the Lagrange multipliers that we have then you can get the true solution the solution of the constrained problem is also in is also a minimum of the penalized problem.

Which means that it suffices to set the penalty parameter to be larger than c star and that solves the problem right. So, this is actually very powerful because it let us effectively get to the solution of the constrained problem through an unconstrained problem and with a finite penalty value penalty parameter value without having to deal with infinities anywhere and.

So, the this is, but the only catch here is that the actual minimization of the penalty of the penal of this function the penalized function the one with the; one with the mod h here in the objective mod h here. The, actual minimization becomes a little problematic because now you have a non differentiable objective the objective is f plus absolute value of h j of x which is not necessarily a differentiable function.

This becomes the this becomes the catch but this is the price you have to pay for having to for getting a strong result of this kind right. So, this is one of the; one of the ways by which you can use the finite value of the penalty and yet use the penalty method to get to the solution of a constrained optimization problem.